



## Pregroups and the French noun phrase

Sylvain Degeilh, Anne Preller

► **To cite this version:**

Sylvain Degeilh, Anne Preller. Pregroups and the French noun phrase. 03023, 2003, pp.19. <lirmm-00191926>

**HAL Id: lirmm-00191926**

**<https://hal-lirmm.ccsd.cnrs.fr/lirmm-00191926>**

Submitted on 26 Nov 2007

**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.

Sylvain Degeilh  
 LIRMM/CNRS  
 Montpellier  
 France  
 degeilh@lirmm.fr

Anne Preller  
 LIRMM/CNRS  
 Montpellier  
 France  
 preller@lirmm.fr

## Pregroups and the French noun phrase

### ABSTRACT:

We study mathematical and algorithmic properties of Lambek's pregroups and illustrate them by analysing the French noun phrase. We establish robustness properties of pregroups and present a simple algorithm of complexity  $n^3$  for deciding reduction in an arbitrary free pregroup as well as a linear algorithm covering a large class of language fragments. In the French noun phrase, the agreement of determiners, adjectives and nouns and word order are treated.

### Introduction

When choosing a formal system for computational purposes, one looks for *efficiency* (taking decidability for granted), *robustness* and *provable verification*. The mathematical concept of pregroups has these properties. As we hope that our readers come from as varied backgrounds as linguistics, logic or computation, we formulate and reformulate mathematical properties in terms of their significance in linguistic analysis, like conservativity of extensions as robustness in Section 1. Section 2 covers algorithmic properties and is more technical. Both Sections are intended as an argument, why one should prefer pregroups to other systems when analysing a language fragment as we did in Section 3 for the French noun phrase.

*Pregroups* are an algebraic tool introduced by J. Lambek [Lamb 99] to recognise grammatically well-formed sentences in natural languages. In linguistic applications, one uses *free pregroups* where the mathematical machinery reduces to two simple schematic rewrite rules, the *contractions*. This tool provides an elegant theory which is universal, independent of the language. Up till now, fragments of Arabic, English, French, German, Italian, Japanese, Latin, Polish, Turkish among others have been analysed with the help of pregroups. Besides the simplicity and the universality, other properties can be proven, due to the mathematical character of the tool.

Indeed, the theory of pregroups is decidable. Therefore, as we shall explain below, the problem whether a given string of words is a well-formed construct of a typed language fragment is decidable. The general decision problem for free pregroups can be implemented by an algorithm of complexity  $n^3$ , guaranteeing thus an efficient algorithm for the typing checking of language fragments. Our "nearest left parentheses" algorithm in Proposition 3 adapts and simplifies an algorithm given in [Earley] for context free grammars. Earley's algorithm was written for context free grammars and uses the number of rewrite rules of the grammar as a bound. It cannot be applied as such to deciding reduction in pregroups. Indeed, the schema of the two contraction rules in pregroups yields an infinite set of rewrite rules in the corresponding context free grammar. Even if only finitely many rewrite rules are needed when type checking a given language fragment,

Earley's algorithm would be quite a bit more cumbersome than ours. In certain language fragments, our algorithm simplifies further to an algorithm which is *linear* in the number of words. We also formulate a sufficient criterion for linearity (Proposition 2), which we hope will influence the way the types are chosen for a language fragment. Finally, typing with pregroups is robust. As extensions of the theory of pregroups are conservative, the typing of a language fragment can be extended to include new constructs, new words without changing the previous analysis. Thus the typing of the French noun phrase given in Section 3 extends earlier work of [Barg-Lamb] on the French sentence structure. It is also part of ongoing work of a more comprehensive fragment. Though our typing covers a liberal notion of determiners and can account for the distinction between the indefinite versus the partitive reading of the articles *du, des* etc., our typing is far from covering all aspects and aims above all at a demonstration of the potential of pregroups.

### 1. Mathematical properties and their linguistic meaning

To explain how pregroups work in linguistic, we recall briefly the definitions and the immediate consequences.

Definition 1: A pregroup is a partially ordered monoid in which each element  $a$  has both a *left adjoint*  $a^\ell$  and a *right adjoint*  $a^r$  satisfying

$$\text{(Contraction)} \quad a^\ell \cdot a \rightarrow 1$$

$$a \cdot a^r \rightarrow 1$$

$$\text{(Expansion)} \quad 1 \rightarrow a^r \cdot a$$

$$1 \rightarrow a \cdot a^\ell .$$

The arrow denotes the partial order relation, the dot denotes multiplication and is generally omitted. By definition, the multiplication is associative and the order is *compatible* with multiplication, that is

$$a \rightarrow b \quad \text{and} \quad c \rightarrow d \quad \text{implies} \quad ac \rightarrow bd .$$

A linguist will work with the *free* pregroup generated by a partially ordered set of so called *basic types*. For a given language fragment, one chooses a set of symbols, the *basic types*, and defines a partial order on this set. We use bold face symbols  $\mathbf{a}, \mathbf{b}$  to vary over basic types.

The basic types and their iterated adjoints form a set  $\Sigma$  of *simple types* :

$$\Sigma = \{ \dots \mathbf{a}^{\ell\ell}, \mathbf{a}^\ell, \mathbf{a}, \mathbf{a}^r, \mathbf{a}^{rr}, \dots \quad \dots, \mathbf{b}^{\ell\ell}, \mathbf{b}^\ell, \mathbf{b}, \mathbf{b}^r, \mathbf{b}^{rr}, \dots \}$$

The strings of simple types will be called *types*. The arrow  $a \rightarrow b$  is now read as “ $a$  reduces to  $b$ ”. In fact, the elements of the free pregroup generated by a set of basic types identify with the types constructed from the same set of basic types, see [Lambek99]. The unit 1 denotes the empty string.

The simple types inherit the order from the basic types as follows:

$$\text{(I)} \quad \text{If } \mathbf{a} \rightarrow \mathbf{b} \text{ then } \mathbf{b}^\ell \rightarrow \mathbf{a}^\ell, \mathbf{b}^r \rightarrow \mathbf{a}^r, \mathbf{a}^{\ell\ell} \rightarrow \mathbf{b}^{\ell\ell}, \mathbf{a}^{rr} \rightarrow \mathbf{b}^{rr}, \mathbf{b}^{\ell\ell\ell} \rightarrow \mathbf{a}^{\ell\ell\ell}, \mathbf{b}^{rrr} \rightarrow \mathbf{a}^{rrr}, \text{ etc.}$$

The order between basic types is “declared”, i.e. can be looked up in a table. Hence the ordering of simple types follows from the same table. Indeed, to check whether  $a \rightarrow b$ , it suffices to check whether  $a$  and  $b$  have the same exponent, i.e.  $a = \mathbf{a}^s$ ,  $b = \mathbf{b}^s$  where  $\mathbf{a}$  and  $\mathbf{b}$  are basic and  $s$  consists of a finite number  $n$  of repetitions of the same symbol, either  $\ell$  or  $r$ . If this is the case, then  $a \rightarrow b$  if either  $n$  is even and  $\mathbf{a} \rightarrow \mathbf{b}$  or  $n$  is odd and  $\mathbf{b} \rightarrow \mathbf{a}$ . Otherwise, neither  $a$  reduces to  $b$  nor  $b$  to  $a$ .

The contractions of simple types can now also be understood as rules :

$$(II) \dots, \mathbf{a}^{\ell\ell} \mathbf{a}^{\ell\ell} \rightarrow 1, \mathbf{a}^{\ell\ell} \mathbf{a}^{\ell} \rightarrow 1, \mathbf{a}^{\ell} \mathbf{a} \rightarrow 1, \mathbf{a} \mathbf{a}^r \rightarrow 1, \mathbf{a}^r \mathbf{a}^{rr} \rightarrow 1, \mathbf{a}^{rr} \mathbf{a}^{rrr} \rightarrow 1, \dots$$

In linguistic applications, where one works with a free pregroup, only the contractions and ordering of simple types is needed (for the detail, see below), i.e. the rules (I) and (II) are all that is to be kept in mind when “typing” a language fragment.

As a first step, the language fragment has to be described in common grammatical terms. We select the words which are supposed to be in the mental or electronic dictionary, for example nouns, adjectives etc. and introduce the grammatical notions describing the grammatically well-formed constructs of the fragment. The next step is to define the set of basic types and its ordering. It must reflect the grammatical constructions under consideration. In general, one will aim at keeping the set of basic types as small and with few inequalities as possible. Finally, to every word of the fragment, one associates one or several types, to be written next to the word in dictionary. Meta-rules serve to organise the content of the dictionary. Instead of explicitly writing the type(s) of a word into the dictionary, it (they) may be described by a meta-rule.

The typing must be done in a way that a sequence of words is a well-formed construct (sentence, noun phrase, etc) if and only if the corresponding string of types reduces to the basic type corresponding to the construct. If a word has more than one type, then it is enough that one of the possible choices yields a string reducing to the type in question. This equivalence is the key property. Every typing that respects it is said to be *correct* (it recognises only well-formed constructs) and *complete* (it recognises all well-formed constructs) with respect to the fragment.

A correct and complete typing satisfies the *principle of Substitution* : If a word is replaced by another word with the same type, then a well-formed string of words remains well-formed. Indeed, being a well-formed string of words is equivalent of having a string of types reducible to a given basic type and, by assumption, the string of types is the same before and after substitution.

Every correct and complete typing satisfies properties which permit to extend the typing to larger fragments without changing the properties of the previous typing. We call these the *robustness* properties:

a) Assigning new types to words, without changing the set of basic types :

Suppose a word  $W$  has type  $d$  and we add a type  $c$  such that  $c \rightarrow d$ . Then every string of words recognised well-formed using the type assignment  $d$  for  $W$  is recognised well-formed using  $c$ .

b) Extensions by new basic types.

This means that one can extend a given set  $\mathbf{B}$  of basic types, by declaring new types and adding inequalities involving the new types, obtaining thus a larger set of basic types  $\mathbf{B}'$ . Then the free pregroup  $\mathbf{P}'$  generated by  $\mathbf{B}'$  includes the free pregroup  $\mathbf{P}$  generated by  $\mathbf{B}$ . Whenever both  $a$  and  $b$  belong to the smaller pregroup  $\mathbf{P}$  and  $a \rightarrow b$  can be derived in  $\mathbf{P}$ , then this also holds in the larger  $\mathbf{P}'$ . *Conservativity* is a sort of converse of this fact, it says : if  $a \rightarrow b$  can be derived in the larger pregroup  $\mathbf{P}'$  and both  $a$  and  $b$  belong to the smaller pregroup, then the whole reduction can be done in the smaller pregroup. This is not trivial, because we might have come to the conclusion  $a \rightarrow b$  by showing  $a \rightarrow c$  and  $c \rightarrow b$  using one of the new types  $c$ . Conservativity is a consequence of the so-called Switching Lemma in [Lambek99] (see below ) provided the order of the old basic types remains unchanged.

The linguistic significance of this is that typing by pregroups is robust. Once it has been shown correct and complete for a fragment it will remain so in the extensions. It also simplifies the task of verifying that a typing is correct and complete. One can proceed step by step, extending the fragment by adding basic type after basic type, and then verify only that the typing involving the new type(s) is also correct and complete for the new constructs.

Here are the promised mathematical details on which our assertions above are based. Except Proposition 1, they are either straight forward or to be found in [Lambek99].

1.  $a1 = a = 1a$  (1 is the unit of the monoid)
2.  $(ab)c = a(bc)$  (multiplication is associative)
3.  $a \rightarrow b$  and  $c \rightarrow d$  implies  $ac \rightarrow bd$  (order is compatible with multiplication).
4.  $ab \rightarrow 1 \rightarrow ba$  implies  $a = b^\ell$  and  $b = a^r$  (adjoints are unique)
5.  $(ab)^\ell = b^\ell a^\ell$ ,  $(ab)^r = b^r a^r$  (adjunction is quasi-distributive)
6.  $a \rightarrow b$  implies  $b^\ell \rightarrow a^\ell$  and  $b^r \rightarrow a^r$  (adjunction reverses the order).
7.  $a^{\ell r} = a = a^{r\ell}$  (no mixed adjoints)

Properties 1. – 3. are in the definition of a monoid, the properties 4. - 7. can be easily derived from Definition 1. For example, to derive that  $b = b^{\ell r}$ , use  $b^\ell b \rightarrow 1 \rightarrow bb^\ell$  and 4. with  $a = b^\ell$ . Similarly, 5. follows from 4. Indeed,  $(ab)(b^r a^r) = a(bb^r) a^r \rightarrow a1a^r = aa^r \rightarrow 1 \rightarrow b^r b = b^r 1b \rightarrow b^r(a^r a)b = (b^r a^r)(ab)$ . Hence, by 4.,  $(ab)^r = b^r a^r$ .

In free pregroups, important additional properties hold. First of all, they are non-commutative and the iterated adjoints of basic types are all different. The most important result is expressed in the Switching Lemma [Lambek99, Proposition 2], which we shall include here for completeness sake. Recall that the elements of the partially ordered set generating the free pregroup are called basic types, their iterated left or iterated right adjoints are called simple types. Strings of simple types are called types, the empty string being denoted 1. Lambek [loc.cit.] shows that this monoid of types is in fact the free pregroup generated by the basic types.

Hence, concatenation plays the role of multiplication and the empty string that of the unit. It is convenient to use a uniform notation for the simple types, be they iterated left adjoints or iterated right adjoints: Write

$$\dots \mathbf{a}^{(-2)}, \mathbf{a}^{(-1)}, \mathbf{a}^{(0)}, \mathbf{a}^{(1)}, \mathbf{a}^{(2)}, \dots$$

for

$$\dots \mathbf{a}^{\ell\ell}, \mathbf{a}^\ell, \mathbf{a}, \mathbf{a}^r, \mathbf{a}^{rr}, \dots$$

Then, by definition, every type has the form

$$\mathbf{a}_1^{(n_1)} \dots \mathbf{a}_k^{(n_k)}$$

where  $\mathbf{a}_1, \dots, \mathbf{a}_k$  are basic types and  $n_1, \dots, n_k$  are integers. The types form already a pregroup, if one defines adjoints by

$$(\mathbf{a}_1^{(n_1)} \dots \mathbf{a}_k^{(n_k)})^\ell := \mathbf{a}_k^{(n_k-1)} \dots \mathbf{a}_1^{(n_1-1)}$$

$$(\mathbf{a}_1^{(n_1)} \dots \mathbf{a}_k^{(n_k)})^r := \mathbf{a}_k^{(n_k+1)} \dots \mathbf{a}_1^{(n_1+1)}$$

and the order on types as the reflexive and transitive closure of the following three relations where  $c, d$  are arbitrary types and  $\mathbf{a}, \mathbf{b}$  are basic:

(Induced step)

$$c \mathbf{a}^{(n)} d \rightarrow c \mathbf{b}^{(n)} d, \text{ if either } n \text{ is even and } \mathbf{a} \rightarrow \mathbf{b} \text{ or } n \text{ is odd and } \mathbf{b} \rightarrow \mathbf{a}.$$

(Generalised contraction)

$$c \mathbf{a}^{(n)} \mathbf{b}^{(n+1)} d \rightarrow cd, \text{ if either } n \text{ is even and } \mathbf{a} \rightarrow \mathbf{b} \text{ or } n \text{ is odd and } \mathbf{b} \rightarrow \mathbf{a}.$$

(Generalised expansion)

$$cd \rightarrow c \mathbf{a}^{(n+1)} \mathbf{b}^{(n)} d, \text{ if either } n \text{ is even and } \mathbf{a} \rightarrow \mathbf{b} \text{ or } n \text{ is odd and } \mathbf{b} \rightarrow \mathbf{a}.$$

Notice that for every basic type  $\mathbf{b}$ , one has the contractions

$$\mathbf{b}^{(n)} \mathbf{b}^{(n+1)} \rightarrow 1, \text{ whatever the integer } n.$$

Switching Lemma (Lambek): Let  $a_1, \dots, a_n$  and  $b_1, \dots, b_m$  be simple types. Then  $a_1 \dots a_n \rightarrow b_1 \dots b_m$  if and only if there are a substring  $a_{i_1} \dots a_{i_k}$  of  $a_1 \dots a_n$  and a substring  $b_{i_1} \dots b_{i_k}$  of  $b_1 \dots b_m$  such that

$$a_1 \dots a_n \rightarrow a_{i_1} \dots a_{i_k} \rightarrow b_{i_1} \dots b_{i_k} \rightarrow b_1 \dots b_m, \quad a_{i_p} \rightarrow b_{i_p}, 1 \leq p \leq k,$$

where  $a_{i_1} \dots a_{i_k}$  is obtained from  $a_1 \dots a_n$  by generalised contractions only,  $b_1 \dots b_m$  is obtained from  $b_{i_1} \dots b_{i_k}$  by generalised expansions only and finally  $b_{i_1} \dots b_{i_k}$  is obtained from  $a_{i_1} \dots a_{i_k}$  by induced steps only.

Call a pair of simple types  $(c, d)$  *contractible*, if  $c = \mathbf{a}^{(n)}$ ,  $d = \mathbf{b}^{(n+1)}$  and either  $\mathbf{a} \rightarrow \mathbf{b}$  and  $n$  even or  $\mathbf{b} \rightarrow \mathbf{a}$  and  $n$  odd.

Corollary (Lambek): For simple types  $a_1, \dots, a_n$  and  $b$ ,  $a_1 \dots a_n \rightarrow b$  holds if and only if there is a simple type  $b' \rightarrow b$  such that one can obtain  $b'$  from  $a_1 \dots a_n$  by repeatedly omitting pairs of contractible adjacent types.

I) The first consequence is the decidability of the *type checking problem*, namely to decide whether  $a \rightarrow b$  where  $a$  is an arbitrary type and  $b$  is simple.

Indeed, given a string  $a_1 \dots a_n$  choose an index  $i$  such that  $(a_i, a_{i+1})$  is contractible and omit  $a_i a_{i+1}$ , start again with  $a_1 \dots a_{i-1} a_{i+2} \dots a_n$  until an irreducible string is reached, i. e. a string without adjacent contractible types. As different choices may lead to different irreducible substrings, the answer is yes, if at least one of the choices leads to a simple type  $b'$  for which  $b' \rightarrow b$ .

Hence, if a correct and complete typing has been provided for a language fragment, then the fragment itself is decidable. Given a string of words in the dictionary of the fragment enumerate all possible strings of types corresponding to the words. As every word has only finitely many types, there are only finitely many such strings. Then test for each string, if it reduces to an appropriate basic type.

Another consequence is the characterisation of conservative extensions :

II) Proposition 1: (Conservativity of extensions) : Let  $\mathbf{B}$  be an ordered subset of  $\mathbf{B}'$ , i.e.  $\mathbf{B}$  is a subset of  $\mathbf{B}'$  and  $a \rightarrow b$  in  $\mathbf{B}$  if and only if  $a, b \in \mathbf{B}$  and  $a \rightarrow b$  in  $\mathbf{B}'$ . Then the free pregroup  $\mathbf{P}'$  generated by  $\mathbf{B}'$  is conservative over the free pregroup  $\mathbf{P}$  generated by  $\mathbf{B}$ . That is to say, for all elements  $e, f$  of  $\mathbf{P}$  such that  $e \rightarrow f$  holds in  $\mathbf{P}'$ ,  $e \rightarrow f$  holds already in  $\mathbf{P}$ .

Proof: Let  $e, f$  be elements of  $\mathbf{P}$  such that  $e \rightarrow f$  holds in  $\mathbf{P}'$ . First consider the special case where  $e \rightarrow f$  is a generalised contraction. Then there are basic types  $a, b \in \mathbf{B}$  and an integer  $n$  such that  $e = c a^{(n)} b^{(n+1)} d$ ,  $f = cd$  and either  $a \rightarrow b$  in  $\mathbf{B}'$  and  $n$  is even or  $b \rightarrow a$  in  $\mathbf{B}'$  and  $n$  is odd. By hypothesis, this implies  $a \rightarrow b$  in  $\mathbf{B}$  and  $n$  is even or  $b \rightarrow a$  in  $\mathbf{B}$  and  $n$  is odd. Therefore  $e \rightarrow f$  is a generalised contraction in  $\mathbf{P}$ . By a similar argument, one shows that a generalised expansion (induced step) in  $\mathbf{P}'$  is also a generalised induced step in  $\mathbf{P}$ .

In the general case, there are simple types  $a_1, \dots, a_n$  and  $b_1, \dots, b_m$  of  $\mathbf{P}$  such that  $e = a_1 \dots a_n$ ,  $f = b_1 \dots b_m$  and  $a_1 \dots a_n \rightarrow b_1 \dots b_m$  in  $\mathbf{P}'$ . Apply the Switching Lemma. There is a substring  $a_{i_1}, \dots, a_{i_k}$  of  $a_1, \dots, a_n$  and a substring  $b_{i_1}, \dots, b_{i_k}$  of  $b_1, \dots, b_m$  such that

$$a_1 \dots a_n \rightarrow a_{i_1} \dots a_{i_k} \rightarrow b_{i_1} \dots b_{i_k} \rightarrow b_1 \dots b_m, \quad a_{i_p} \rightarrow b_{i_p} \text{ in } \mathbf{P}', \quad 1 \leq p \leq k,$$

$a_{i_1} \dots a_{i_k}$  is obtained from  $a_1 \dots a_n$  by repeated generalised contractions and  $b_1 \dots b_m$  is obtained from  $b_{i_1} \dots b_{i_k}$  by repeated generalised expansions. Now for every generalised contraction  $c \rightarrow d$ , from  $c \in \mathbf{P}$  follows  $d \in \mathbf{P}$ . Indeed,  $d$  is obtained from  $c$  by omitting two simple types. Similarly, for every generalised expansion  $c \rightarrow d$ , from  $d \in \mathbf{P}$  follows  $c \in \mathbf{P}$ . Hence  $a_{i_1} \dots a_{i_k}$  and  $b_{i_1} \dots b_{i_k}$  are in  $\mathbf{P}$  and all the intermediary generalised

contractions, expansions and induced steps take place in  $\mathbf{P}$ . By the first part of the proof, this implies that  $a_1 \dots a_n \rightarrow a_{i_1} \dots a_{i_k} \rightarrow b_{i_1} \dots b_{i_k} \rightarrow b_1 \dots b_m$ , and  $a_{i_p} \rightarrow b_{i_p}$ ,  $1 \leq p \leq k$ , already hold in  $\mathbf{P}$ .

The Proposition 1 simplifies the task of verifying that the typing is correct and complete. One can proceed step by step, extend the fragment by adding new basic types and assigning new types to words. Then to show that the typing is correct and complete with respect to the larger fragment, it suffices to consider the sequences of words with new types. The only prerequisite is not to change the order between the old basic types. For example, suppose  $\mathbf{B}' = \mathbf{B}[c]$  where  $c$  is a new basic type. We may declare  $a \rightarrow c$  and/or  $c \rightarrow b$  for some “old” basic types  $a, b \in \mathbf{B}$ . However, we must take care not to declare both  $a \rightarrow c$  and  $c \rightarrow b$ , unless  $a \rightarrow b$  already holds in  $\mathbf{B}$ . Notice that this step by step approach is normally taken for granted. The conservativity property says that this approach is safe. This is not as trivial as one might think. Obviously, a sequence of words which is a grammatically well-formed construct of the smaller fragment remains so in the larger fragment. Just use the typing and the reduction in the smaller pregroup  $\mathbf{P}$ , which remains a reduction in the bigger one  $\mathbf{P}'$ . It is, however, not so obvious that a sequence of words which is not well-formed in the smaller fragment does not become well-formed in the bigger fragment, even if no new types are involved. Indeed, suppose that a sequence of words gets assigned a type  $c$  in the smaller fragment, but is not well-formed. As the typing in  $\mathbf{P}$  is correct and complete with respect to the smaller fragment, we have  $c \not\rightarrow a$  where  $a$  is the basic type corresponding to the grammatical notion under investigation. A priori, one cannot exclude that  $c \rightarrow d$  and  $d \rightarrow a$ , where  $d$  is a new type in  $\mathbf{P}'$ . But then, it would follow that  $c \rightarrow a$  in  $\mathbf{P}'$ . From this we would have to conclude that the sequence of words is now well-formed, because its type reduces to  $a$ . Conservativity guarantees that this cannot happen, as from  $c \rightarrow a$  in  $\mathbf{P}'$  it would follow that  $c \rightarrow a$  in  $\mathbf{P}$ , contradicting  $c \not\rightarrow a$ .

This property is used in fact continuously and most of the time without saying so. We will do so in Section 3

## 2. Algorithmic Properties

The type checking problem for free pregroups is to decide by a general method whether  $a_1 \dots a_n \rightarrow b$  for arbitrary simple types  $a_1, \dots, a_n$  and any irreducible type  $b$ , i.e. a string of simple types containing no contractible pair of adjacent simple types. The Corollary to the Switching Lemma gives such a method:

For deciding whether  $a_1 \dots a_n \rightarrow b$  holds, it suffices to omit a contractible pair of adjacent types  $a_i a_{i+1}$  and repeat this procedure in all possible ways until the resulting string is irreducible. If at least one of the resulting strings  $b'$  satisfies  $b' \rightarrow b$ , the answer is yes, otherwise it is no.

Any implementation of this decision procedure is called a *type checking* algorithm. The decision procedure for free pregroups provides a solution to the problem whether a sequence of words belongs to a language fragment. For each word in the sequence choose one of its types in the dictionary. The sequence of words belongs to the fragment if and only if one of the strings of assigned types reduces to a specific basic type, say the type of a noun phrase or the type of a sentence. An algorithm which provides associated strings of types for a string of words is called a *type assignment* algorithm. To keep this paper in reasonable limits, we restrict ourselves to the efficiency of type checking algorithms.



The result of the reduction of a string of types to an irreducible one depends in general on the choice of the contracted pairs. For example,  $a^\ell a^\ell a^r$  reduces with two contractions to the empty string, namely by contracting the left most and the right most pair. On the other hand, by contracting the central pair  $a^\ell a$  one obtains the irreducible string  $a^\ell a^r$ . The type checking problem for free pregroups can be replaced by the more specific one which checks if a string reduces to the empty string 1, see Lemma 1 below. The set of all strings  $a_1 \dots a_n$  such that  $a_1 \dots a_n \rightarrow 1$  can be characterised by a context free grammar. However, this context free grammar has an infinite set of terminal symbols (the simple types and 1) and its set of rewrite rules is infinite. Indeed, it includes all the rules of the form  $E \Rightarrow E a^{(n)} E a^{(n+1)} E$  where  $a$  is a basic type,  $n$  an integer and  $E$  a new, non-terminal symbol standing for a string which reduces to the empty one. Hence the efficient algorithm of complexity  $K n^3$  in [Earley] does not apply, because  $K$  is a common bound for the set of terminal symbols and the set of rewrite rules. If we are only interested in linguistic applications, we may restrict the set of strings by requiring that the intervening simple types belong to a given finite set. However, our algorithm in Proposition 3 not only is simpler than Earley's but also implements the decision procedure for type checking in an arbitrary free pregroup. Like Earley's, it reads the string from left to right and has complexity  $n^3$ , but with a much lower constant. For certain sets of strings, our algorithm can be simplified even further to one that is linear in the length  $n$  of the string, see Proposition 4. This linear algorithm is useful in linguistic applications. The strings of types associated to a word of the language fragment very often satisfy the criterion which permits the use of the linear algorithm.

We begin with a criterion for subsets of free pregroups permitting linear type checking. A triple of simple types  $(a, b, c)$  is said to be *critical*, if both  $(a, b)$  and  $(b, c)$  are contractible. This is equivalent to saying that there are basic types  $a, b, c$  and an integer  $n$  such that  $a = a^{(n-1)}$ ,  $b = b^{(n)}$ ,  $c = c^{(n+1)}$  and  $b^{(n)} \rightarrow a^{(n)}$ ,  $b^{(n)} \rightarrow c^{(n)}$ . The latter is the case, exactly when either  $n$  is even and  $b \rightarrow a, b \rightarrow c$  or  $n$  is odd and  $a \rightarrow b, c \rightarrow b$ . A string of simple types is *linear*, if it has no substring of length 3 which is a critical triple. The typing of a language fragment is said to be *linear*, if all strings of types corresponding to strings of words in the dictionary are linear. For example,  $a^r a a^\ell$  is linear, but  $a^\ell a a^r$  is not linear. Or if  $a \rightarrow b$ , then  $c b a^\ell a a b^r$  is linear, but if  $a \rightarrow b$ , then  $c b a^\ell a a b^r$  is not linear.

Proposition 2: Every linear type has a unique irreducible form.

Proof: Suppose  $a_1 \dots a_n$  is linear. Use induction on  $n$ , the length of the string. If  $n = 1$ , the property is obvious. For the induction step, notice that every substring of  $a_1 \dots a_n$  is again linear. Moreover, whenever  $a_i a_{i+1}$  and  $a_j a_{j+1}$  are both contractible, the indices  $i, i+1, j, j+1$  are all different. Suppose  $a_1 \dots a_n$  has  $k$  pairs of contractible adjacent types. Omitting them in  $a_1 \dots a_n$  corresponds to  $k$  simultaneous contractions. They can be done in any order without changing the result. If  $k = 0$ , then  $a_1 \dots a_n$  is irreducible. Otherwise, the unique substring resulting from the  $k$  contractions has length less than  $n$  and we may conclude by induction hypothesis.

Lemma 1: Suppose that  $a_1, \dots, a_n$  and  $b$  are simple types. Then  $a_1 \dots a_n \rightarrow b$  if and only if  $b^\ell a_1 \dots a_n \rightarrow 1$ .

Proof: Assume that  $b^\ell a_1 \dots a_n \rightarrow 1$  and use induction on the length  $n$ . If  $n=1$ ,  $b^\ell a_1 \rightarrow 1$  must be generalised contraction. From the definition of generalised contractions, it follows at once that  $a_1 \rightarrow b$ . Assume  $n > 1$ . If  $b^\ell a_1$  are contractible, we conclude as above. Otherwise, the first generalised contraction transforming  $b^\ell a_1 \dots a_n$  to 1 must be in  $a_1 \dots a_n$ . Hence the result has the form  $b^\ell a_{i_1} \dots a_{i_{n-2}}$  with  $b^\ell a_{i_1} \dots a_{i_{n-2}} \rightarrow 1$ . By induction hypothesis,  $a_{i_1} \dots a_{i_{n-2}} \rightarrow b$ . As  $a_1 \dots a_n \rightarrow a_{i_1} \dots a_{i_{n-2}}$ , the property follows. As  $a_1 \dots a_n \rightarrow b$  implies  $b^\ell a_1 \dots a_n \rightarrow b^\ell b \rightarrow 1$ , we are done.

It is easy to device a linear-time algorithm which reads a string from left to right and produces an irreducible form. For linear strings, it will provide the only irreducible string obtainable by contractions. In fact, this left to right algorithm is a closely linked to the more general one below which solves the type checking problem without restriction on the strings.

Proposition 3 : There is an algorithm which decides in at most  $n^3$  steps whether  $a_1 \dots a_n \rightarrow 1$ .

Proof: We define a function  $Nlp$  which maps a string  $S$  of length  $n$  and an integer  $i \in \{1, \dots, n+1\}$  to a subset of  $\{1, \dots, i\}$  and show that

- 1)  $S \rightarrow 1$  if and only if  $0 \in Nlp(S, n+1)$
- 2)  $Nlp(S, n+1)$  can be calculated in at most  $n^3$  steps

We omit the parameter  $S$  to simplify notation. The intuitive idea of  $Nlp$  is the following : Process the string  $S = a_1 \dots a_n$  from left to right, looking for all possible contractions, reading the symbol  $a_i$  at stage  $i+1$ . This symbol could be a left “parenthesis” to a later  $a_p$ , i.e.  $a_i \leq a_p^\ell$ , or a “right parenthesis” to some earlier  $a_j$ , i.e.  $a_i \leq a_j^r$ . As it cannot be both in the same reduction, one has to keep track of more than one reduction, or at least of the indispensable information. It turns out that the left parentheses ready for contraction at stage  $i$  is all that is needed to continue the processing to stage  $i+1$ . Therefore at stage  $i+1$ ,  $a_i$  is declared an “open” left parenthesis (to be contracted with some right parenthesis which might come up later), i.e. the index  $i$  is stored in  $Nlp(i+1)$ . Notice that for this choice,  $a_i$  becomes the nearest open left parenthesis for  $a_{i+1}$ . However,  $a_i$  might also be a right parenthesis to the nearest left parenthesis in a reduction made up to stage  $i$ . So, for each  $j \in Nlp(i)$ , we check if  $a_i \leq a_j^r$ . If this is the case,  $a_j$  becomes “closed” and the open left parenthesis at stage  $j$  become ready again for contraction at stage  $i+1$ . Hence,  $Nlp$  has the following definition.

Definition 2:

$$Nlp(1) = \{0\}$$

$$Nlp(i+1) = \{i\} \cup \bigcup_{j \in Nlp(i), a_j \leq a_i^r} Nlp(j), \quad 1 \leq i \leq n$$

The rest of the proof follows from the next 5 lemmas.

**Lemma 2 :**  $Nlp(i+1) \subseteq \{0, \dots, i\}$ .

This follows immediately from the definition.

**Lemma 3 :** If  $k \in Nlp(i+1)$ , then  $a_{k+1} \dots a_i \rightarrow I$ , i.e.  $a_k$  is ready for contraction with  $a_{i+1}$ .

Indeed, use induction on  $i-k$ . If  $i=k$ , then  $a_{k+1} \dots a_i$  is the empty string and  $I \rightarrow I$ . Assume  $i-k > 0$  and  $k \in Nlp(i)$ . Then there is a  $j \in Nlp(i)$  such that  $k \in Nlp(j)$  and  $a_i \leq a_j^r$ . By Lemma 2 and assumption,  $k < j < i$ , therefore  $j-k < i-k$  and  $i-j < i-k$ . Thus by induction hypothesis,  $a_{k+1} \dots a_{j-1} \rightarrow I$  and  $a_{j+1} \dots a_{i-1} \rightarrow I$ . Finally,  $a_{k+1} \dots a_{j-1} a_j a_{j+1} \dots a_{i-1} a_i \rightarrow a_j a_i \rightarrow a_j a_j^r \rightarrow I$ .

**Lemma 4 :** If  $a_{k+1} \dots a_i \rightarrow I$ , then  $k \in Nlp(i+1)$ .

Again proceed by induction on  $i-k$ . If  $i=k$ , then  $k \in Nlp(i+1)$  by definition. Let  $i-k > 0$  and assume that  $a_{k+1} \dots a_i \rightarrow I$ . Then there is an index  $j$  such that  $k+1 \leq j \leq i-1$  such that  $a_{j+1} \dots a_{i-1} \rightarrow I$ ,  $a_i \leq a_j^r$  and  $a_{k+1} \dots a_{j-1} \rightarrow I$ . By induction hypothesis, this implies that  $j \in Nlp(i)$  and  $k \in Nlp(j)$ . So,  $k \in Nlp(i+1)$  by definition.

**Lemma 5 :**  $0 \in Nlp(n+1)$  if and only if  $a_1 \dots a_n \rightarrow I$ .

Apply Lemma 3 and Lemma 4 to  $k=0$ ,  $i=n$ .

**Lemma 6 :**  $Nlp(i+1)$  can be calculate from  $Nlp(i)$  in at most  $i^2 + 1$  steps. The complexity of  $Nlp(n+1)$  is  $n^3$ .

Indeed, to calculate  $Nlp(i+1)$  we first copy  $i$  into  $Nlp(i+1)$  and then must compare  $a_i$  with  $a_j^r$  for every  $j \in Nlp(i)$ . By Lemma 2, we make at most  $i$  comparisons. Each time the comparison succeeds, the elements of  $Nlp(j)$  are copied into  $Nlp(i+1)$ . But  $Nlp(j)$  has at most  $j$  elements and  $j$  is bounded by  $i$ . On the whole, we have executed at most  $i^2 + 1$  steps, counting a comparison and recopying an element as one-step operation. Finally, to calculate  $Nlp(n+1)$ , we must calculate  $Nlp(1), \dots, Nlp(n)$  one after the other, so we do at most  $n(n^2 + 1)$  steps. Notice that the constant of  $n^3$  is 1 in this estimate.

**Corollary 1:** If there is a bound  $K$  such that  $Nlp(i) < K$  for all  $i$ , then the complexity of  $Nlp(n+1)$  is linear  $n$ .

Indeed, the number of steps from  $Nlp(i)$  to  $Nlp(i+1)$  is bounded by  $K^2 + 1$ , hence  $Nlp(n+1)$  can be calculated in at most  $(K^2 + 1)n$  steps.

Notice that  $Nlp$  makes enough calculations to scan all possible reductions of the string. In some situations, it might be enough to find just one reduction yielding an irreducible string. The hope is that if there are fewer reductions to keep track of, it is more likely to find a bound for the set of “nearest left open parentheses”. A slight modification of  $Nlp$  produces such a reduction, which can be calculated in at most  $4n$  steps. The idea is to open a new left parenthesis, only if necessary. Call the corresponding function  $Llp$ , the “lazy left parenthesis” function. It is defined on  $\{1, \dots, n\}$ , the string  $S$  of length  $n$  being given, and takes subsets of  $\{1, \dots, n\}$  as values.

Definition 3:

$$Llp(1) = \{0\}$$

$$Llp(i+1) = \begin{cases} \bigcup_{j \in Link(i)} Llp(j), & \text{if } Link(i) \neq \emptyset \\ \{i\} & \text{else} \end{cases}$$

where  $Link(i) = \{j \in Llp(i) : 1 \leq j, a_i \leq a_j^r\}$ ,  $1 \leq i \leq n$ .

One shows easily that  $Llp(i+1)$  is included in  $Nlp(i+1)$  and has exactly one element. Moreover,  $j \in Link(i)$  means that  $a_j$  is contracted with  $a_i$  in the reduction defined by  $Llp$ . More precisely, say that  $i$  and  $j$  are *linked*, if  $j \in Link(i)$  or  $i \in Link(j)$ . Then :

Lemma 7 : 1) If  $k \in Llp(i+1)$ , then  $a_{k+1} \dots a_i \rightarrow 1$  and whenever  $k+1 \leq m \leq i$  there exists a  $p$  such that  $m$  and  $p$  are linked and  $k+1 \leq p \leq i$ .

2) If  $j \in Link(i)$ , then  $a_j \dots a_i \rightarrow 1$ .

3)  $m \in Link(p)$  implies for all  $i > p$  and all  $j \in Llp(i)$  that  $j > p$  or  $j < m$ .

4) Every index  $p$  is linked to at most one index  $m$ .

Proof: 1) Proceed as in Lemma 3. Assume  $k \in Llp(i+1)$  and  $k+1 \leq m \leq i$ . If  $k = i$ , there is nothing to show. If  $k < i$ , then there is  $j \in Link(i)$  such that  $k \in Llp(j)$ . As  $Link(i) \subseteq Llp(i)$ , the induction hypothesis applies to  $i-j$  and to  $j-k$ . Remark that if  $m = j$  or  $m = i$ , then  $m$  is indeed linked to some  $p$  with  $k+1 \leq p \leq i$ . 2) is an immediate consequence of 1). 3) Suppose  $m \in Link(p)$ ,  $i > p$  and  $j \in Llp(i)$ . Use induction on  $i-p$ . In the case of  $i = p+1$ , we have  $j \in Llp(m)$ , as  $Link(p) \neq \emptyset$ . Hence  $j < m$ . If  $i > p+1$ , then either  $j = i-1 > p$  or  $j \in Llp(i-1)$  and therefore by induction hypothesis  $j > p$  or  $j < m$ . To see 4), remark first that an index  $p$  cannot be linked to two different smaller indices, as  $Link(p)$  has at most one element. By 3)

an index  $m$  cannot be linked to two larger ones, say  $p$  and  $i$  with  $p < i$ . And it also cannot be linked to a smaller and a larger one, as  $m \in Link(p)$ ,  $p \in Link(i)$  would also contradict 3).

Finally, the next and last Lemma confirms that the unlinked elements of the string in increasing order form an irreducible substring.

Lemma 8 : Let  $Unlinked = \{i_1, \dots, i_q\}$  be the set of unlinked indices in increasing order. Then the following holds

I) every index less than  $i_1$  (respectively between  $i_\ell$  and  $i_{\ell+1}$ , respectively larger than  $i_q$ ) is linked to some index below  $i_1$  (respectively between  $i_\ell$  and  $i_{\ell+1}$ , respectively larger than  $i_q$ ). II)  $a_{i_\ell}$  and  $a_{i_{\ell+1}}$  are not contractible.

Proof: The assertion I) follows from Lemma 7, 1) and 2). To show II), assume  $i = i_{\ell+1}$  and let  $j \in Llp(i)$ . By choice of  $i$ ,  $a_i \not\leq a_j^r$ , i.e.  $a_j$  and  $a_i$  are not contractible. Hence, it suffices to show that  $j = i_\ell$ . In view of Lemma 7, 1), it suffices to show that  $j$  is not linked to any index. If  $j$  was linked to a smaller index, this would contradict Lemma 7, 3). If  $j$  was linked to an index greater than  $i$ , this would imply that  $i$  is linked by Lemma 7, 1). Finally,  $j$  cannot be linked to a larger index which would be less than  $i$ , because of Lemma 7, 1) and 4).

Proposition 4 : The type checking problem of a linear fragment can be decided by a linear algorithm.

Proof: By Lemma 1, it suffices to find an algorithm which is linear in the length of the string and produces an irreducible form of the string. For this it suffices to add a new step at stage  $i+1$ : erase  $j$  and  $i$  from  $\{1, \dots, n\}$ , whenever the test  $j \in Link(i)$  succeeds. Together with the two steps to calculate  $Llp(i+1)$ , at most four operations are performed on the whole at stage  $i+1$ .

### 3. French Noun phrases

By Proposition 1 on conservativity of extensions, the typing proposed below may be viewed as an extension of the typing in [Barg-Lamb]. To make things work, the new basic types have to be related to the old ones. For example, we introduce the basic type  $n_{gn}$  to denote a complete noun phrase, depending on its *gender*  $g$  and *number*  $n$ . We postulate  $n_{gn} \rightarrow n$ , where  $n$  is a type used in [loc.cit.] in the situations where the gender and number are of no importance.

The noun phrases analysed here are either names or a determiner followed by a noun. After the noun or between it and the determiner may be adjectives. A noun alone is in general an incomplete noun phrase, though French knows some noun phrases formed from a noun without a determiner. As these are rather archaic and exceptional, we will not consider them in what follows. The so-called prenominal adjectives precede the noun, the postnominal adjectives follow it. A noun with correctly declined and correctly placed adjectives forms an incomplete noun phrase. A determiner transforms an incomplete noun phrase into a complete noun phrase, which may be subject or object in a sentence. The notion of determiner follows [Le bon usage], it includes the

indefinite article *un, une*, the definitive article, *le, la, l', les* and its contracted forms with *de* namely *du, des*, the possessive and demonstrative pronouns such as well as the preposition *de* followed by an adverb of degree like *beaucoup, peu*

### 3.1 Nouns and Adjectives

Nouns are count-nouns like *chat, pomme* and mass-nouns like *eau, pain, vent* etc, but also *courage, beauté* and so on. Noun phrases inherit this distinction. Nouns as well as adjectives vary in gender and number. Therefore we distinguish the types of noun phrases by indices  $g$  and  $n$ , where  $g = 1$  stands for masculine  $g = 2$  for feminine, whereas  $n = 1$  means singular and  $n = 2$  means plural. Thus we use the symbol  $c_{gn}$  for a count-noun phrase and similarly  $m_{gn}$  for a mass-noun phrase. The gender of a noun is given in the dictionary as well as its plural form. Number and gender of noun and adjectives must agree in an noun phrase and constitute the number and gender of the noun phrase. Many mass-nouns have no plural form, for example *riz*. If it has one, we treat the plural form as count-noun.

The formation of the noun phrase also has to take into account the dislike of French speakers of the hiatus where two consecutive vowels would clash. For example, *l'arbre, cet arbre, bel arbre, mon eau*, but *\*le arbre, ce arbre, \*beau arbre, \*ma eau*.<sup>1</sup> Every type  $x$  will have a copy  $x'$  used for words or sequences of words of which the first letter is a vowel.<sup>2</sup> Thus we have basic types  $x_{gn}$ , where  $x$  stands for  $c, c', m$  or  $m'$ .

Adjectives are not only divided into *prenominal* adjectives and *postnominal* adjectives<sup>3</sup>, but they also must respect a certain order, if more than one precedes or follows the noun. For example *bon vin, vin blanc, beau petit chat, autre beau petit chat, vin blanc pétillant, bon petit vin blanc pétillant* etc. but *\*vin bon, \*blanc vin, \*petit beau chat, \*vin pétillant blanc* etc. If two adjectives should occupy the same position, they must be linked by a copula, for example *film noir et blanc*.<sup>4</sup> We assume that the classification of the adjectives into prenominal and postnominal hierarchy classes is known and can be looked up in the dictionary. We use Arabic digits for the prenominal classes  $C_1, C_2, \dots$ , Roman ones for the postnominal classes  $C_I, C_{II}, \dots$ , i.e. we have classes  $C_h$  where  $h \in \{1, 2, \dots\} \cup \{I, II, \dots\}$ . The lower the number of its class, the closer the adjective will be to the noun. Thus *blanc* belongs to  $C_I$ , *petit* to  $C_I$ , *pétillant* to  $C_{II}$ . This distribution into hierarchy classes necessitates the introduction of corresponding basic types for the incomplete noun phrases  $x_{hgn}$  where  $h \in \{1, 2, \dots\} \cup \{I, II, \dots\}$  or  $h = 0$  for bare nouns:

We postulate

$$x_{hgn} \rightarrow x_{gn}, \text{ for all } h, x = c, c', m \text{ or } m'.$$

<sup>1</sup> This phenomenon is even more pervasive in the spoken language where the otherwise silent terminal consonant of a word is pronounced, if the following word starts with a vowel.

<sup>2</sup> There are words whose first letter is a silent  $h$  which are assimilated to words starting with a vowel.

<sup>3</sup> Some adjectives may belong to both classes, especially if classic French or regional variations are also to be covered

<sup>4</sup> This can be done with the usual polymorphic typing of the copula. To keep the paper in reasonable limits, we ignore this case.

The type  $x_{0gn}$  corresponds to a noun, it is given in the dictionary, for example

<i>amande</i> :	$\mathbf{c}'_{021}$
<i>pain</i> :	$\mathbf{m}_{011}$
<i>eau</i> :	$\mathbf{m}'_{021}$

The types  $x_{ign}$ , for  $i \in \{1, 2, \dots\}$  correspond to incomplete noun phrases starting with an adjective in hierarchy class  $C_i$ . The types  $x_{kgn}$  for  $k \in \{I, II, \dots\}$  correspond to incomplete noun phrases starting with a noun and ending with an adjective in class  $C_k$ . This will be accomplished by assigning the types to adjectives as follows:

Meta-Rule 1 (Special adjectives): The masculine form singular of a few special adjectives like *beau*, *nouveau*, *vieux* has a variant *bel*, *nouvel*, *viel* to be used if the next word starts with a vowel.<sup>5</sup> The type of *beau*, *bel* etc. is therefore

<i>beau</i> , ...	: $\mathbf{x}_{211}\mathbf{x}_{h11}^\ell$ , for $h = 1, 0, I, II, \dots$
<i>vieux</i> , ...	: $\mathbf{x}_{111}\mathbf{x}_{h11}^\ell$ , for $h = 0, I, II, \dots$
<i>bel</i> , ...	: $\mathbf{x}_{211}\mathbf{x}'_{h11}^\ell$ , $h = 1, 0, I, II, \dots$
<i>viel</i> ,	: $\mathbf{x}_{111}\mathbf{x}'_{h11}^\ell$ , for $h = 0, I, II, \dots$

where  $\mathbf{x} = \mathbf{c}$ ,  $\mathbf{x} = \mathbf{m}$ .

The types of the non special adjectives and the feminine singular and the plural form of the special adjectives are described by the meta-rule below:

---

<sup>5</sup> In the spoken language, every adjective ending in a silent consonant will have a variant form where the last letter is audible if the following word starts with a vowel.

Meta-Rule 2 (Adjectives):

Let  $A$  be an adjective and  $A_{gn}$  its declined form of gender  $g$  and number  $n$ .

1) If  $A$  belongs to the prenominal hierarchy class  $C_i$ ,  $i = 1, 2, \dots$ , then

$$A_{gn} : \mathbf{x}_{i gn} \mathbf{y}_{h gn}^{\ell}, \quad h = i-1, \dots, 0, I, II, \dots,$$

where

if  $A$  starts with a consonant, either  $\mathbf{x} = \mathbf{c}$  and  $\mathbf{y} = \mathbf{c}, \mathbf{c}'$  or  $\mathbf{x} = \mathbf{m}$  and  $\mathbf{y} = \mathbf{m}, \mathbf{m}'$ ,

if  $A$  starts with a vowel, either  $\mathbf{x} = \mathbf{c}'$ ,  $\mathbf{y} = \mathbf{c}, \mathbf{c}'$  or  $\mathbf{x} = \mathbf{m}'$ ,  $\mathbf{y} = \mathbf{m}, \mathbf{m}'$ .

2) If  $A$  belongs to the postnominal hierarchy class  $C_k$ ,  $k = I, II, \dots$ , then

$$A_{gn} : \mathbf{x}_{h gn}^r \mathbf{x}_{i gn}, \quad \text{where } \mathbf{x} = \mathbf{c}, \mathbf{c}', \mathbf{m}, \mathbf{m}', \quad 0 \leq h < k.$$

Examples:

$$\text{vin} : \mathbf{m}_{011}$$

$$\text{blanc} : \mathbf{c}_{011}^r \mathbf{c}_{I11}, \mathbf{m}_{011}^r \mathbf{m}_{I11}$$

$$\text{pétillant} : \mathbf{c}_{011}^r \mathbf{c}_{II11}, \mathbf{c}_{I11}^r \mathbf{c}_{II11}, \mathbf{m}_{011}^r \mathbf{m}_{II11}, \mathbf{m}_{I11}^r \mathbf{m}_{II11},$$

$$\text{bon} : \mathbf{c}_{211} \mathbf{c}_{011}^{\ell}, \mathbf{c}_{211} \mathbf{c}_{I11}^{\ell}, \mathbf{c}_{211} \mathbf{c}_{II11}^{\ell}, \dots, \mathbf{m}_{211} \mathbf{m}_{011}^{\ell}, \mathbf{m}_{211} \mathbf{m}_{I11}^{\ell}, \mathbf{m}_{211} \mathbf{m}_{II11}^{\ell}, \dots$$

*vin blanc*

$$\mathbf{m}_{011} \mathbf{m}_{011}^r \mathbf{m}_{I11} \rightarrow \mathbf{m}_{I11}$$

*vin pétillant*

$$\mathbf{m}_{011} \mathbf{m}_{011}^r \mathbf{m}_{II11} \rightarrow \mathbf{m}_{II11}$$

*vin blanc pétillant*

$$\mathbf{m}_{011} \mathbf{m}_{011}^r \mathbf{m}_{I11} \mathbf{m}_{I11}^r \mathbf{m}_{II11} \rightarrow \mathbf{m}_{II11}$$

*bon vin blanc pétillant*

$$\mathbf{m}_{211} \mathbf{m}_{II11}^{\ell} \mathbf{m}_{011} \mathbf{m}_{011}^r \mathbf{m}_{I11} \mathbf{m}_{I11}^r \mathbf{m}_{II11} \rightarrow \mathbf{m}_{211} \mathbf{m}_{II11}^{\ell} \mathbf{m}_{II11} \rightarrow \mathbf{m}_{211}.$$

Similarly,

$$\text{petit chat, petit chat noir: } \mathbf{c}_{111}$$

$$\text{beau (petit) chat (noir): } \mathbf{c}_{211}$$

A comment on our use of indices is appropriate: Notice that the types of *\*petite chat* or *\*chat petit* will not reduce to a simple type. Types which differ “only” by the value of an index, say  $\mathbf{c}_{011}$  and  $\mathbf{c}_{111}$ , are just as



different types as those which look “completely different”, say  $m_{III}$  and  $c_{012}$ . The use of indices is convenient, when defining the dictionary, i.e. when assigning types to the words of the language fragment. Each index, better each position of the index in the subscript of a type symbol represents a “feature” of the concept, like gender, number, position of the adjective. The fact that a (sequence of) word(s) starts with a vowel or consonant could also reasonably be called a feature. For reasons of readability, this feature is expressed by the presence or absence of the symbol ' in the superscript, (prime). The theory of pregroups does not include unification of features, but indices are a good device to find efficient type assignment algorithms .

When checking whether *bon vin* is an incomplete noun phrase, we must try all possible type assignments of this sequence of two words until we hit one which reduces to the appropriate basic type. Only one of the possible types for *bon*, namely  $m_{211}m_{011}^\ell$ , will result in a string which reduces the type of *bon vin* to  $m_{211}$ . Notice that another type for *bon* must be used in *bon vin blanc pétillant*. Type assignment can be made more efficient by keeping the variables as long as possible. The Meta-rule describes types for *bon* as  $x_{211}x_{h11}^\ell$ , for  $h=1,0,I,II,\dots$ ,  $x=c,m$ . This corresponds to eight or more types, depending on the number of hierarchy classes. As a first step, the improved type assignment algorithm would assign the string  $x_{211}x_{h11}^\ell m_{011}$  to *bon vin* and at a second step make  $x=m$  and  $h=0$ . Hence, a simple calculus of equality of “features” is part of an efficient type assignment algorithm.

### 3.2. Determiners

A determiner transforms an incomplete noun phrase into a complete noun phrase, which may be subject or object in a sentence. The notion of determiner follows [Le bon usage], it includes the indefinite article *un, une*, the definitive article, *le, la, l', les* and its contracted forms with *de* namely *du, des*, the possessive and demonstrative pronouns such as well as the preposition *de* followed by an adverb of degree like *beaucoup, peu*

The complete noun phrase formed with the possessive or demonstrative pronoun, definite or indefinite article can be subject or attribute or direct object. It has type  $n_{gn}$  or  $n'_{gn}$  where the indices  $g, n$  represent gender and number. If the latter do not matter, we use  $n$  with  $n_{gn} \rightarrow n$ . In view of Proposition 1 our work can be considered as an extension of the analysis given in [Barg-Lamb]. In some constructions with the preposition *de*, one needs to know if the complete noun phrase is formed with a mass-noun or a count-noun. Therefore it is convenient to introduce a type  $\bar{y}_{gn} \rightarrow n_{gn}$ , for  $y=c, m$  and  $\bar{y}_{gn} \rightarrow n'_{gn}$ , for  $y=c', m'$ .

#### General complete noun phrase

Roughly speaking, names are complete noun phrases and so are nouns, with or without adjectives, when preceded by an article or a demonstrative or possessive pronoun.

*Albert* :  $\bar{c}'_{11}$

*Marie* :  $\bar{c}_{21}$

<i>le</i>	:	$n_{11}x_{11}^{\ell}$ , $x = c, m$
<i>les</i>	:	$n_{g2}x_{g2}$ , $x = c$ , $g = 1, 2$
<i>ces, mes, ...</i>	:	$\bar{x}_{g2}x_{g2}$ , $x = c$ , $g = 1, 2$
<i>ce, mon, ton, son, notre, votre, leur</i>	:	$\bar{x}_{11}x_{11}^{\ell}$ , $x = c, m$
<i>la, cette, ma, ta, sa, notre, votre, leur</i>	:	$\bar{x}_{21}x_{21}^{\ell}$ , $x = c, m$
<i>cette, mon, ton, son, notre, votre, leur</i>	:	$\bar{x}_{21}x'_{21}^{\ell}$ , $x = c, m$
<i>l'</i>	:	$\bar{x}_{g1}x'_{g1}^{\ell}$ , $g = 1, 2$ , $x = c, m$
<i>cet</i>	:	$\bar{x}_{11}x'_{11}^{\ell}$ , $x = c, m$
<i>ces, mes, ...</i>	:	$\bar{x}_{g2}x_{g2}$ , $x = c$ , $g = 1, 2$
<i>un</i>	:	$\bar{x}'_{11}x_{11}^{\ell}$ , $x = c, m$
<i>une</i>	:	$\bar{x}'_{21}x_{21}^{\ell}$ , $x = c, m$

The difference between the types of *le, les* and the other determiners lies in the fact that prepositions like *de, à* contract with *le, les* to yield a new word: *du*, (*\*de le*), *des* (*\*de les*) etc.

Example :

$$\begin{array}{c}
 \text{un} \quad \text{bon} \quad \text{vin} \quad \text{blanc} \\
 \bar{m}'_{11}m_{11}^{\ell} \quad m_{211}m_{111}^{\ell} \quad m_{011} \quad m_{011}^r m_{111} \quad \rightarrow \quad \bar{m}'_{11}m_{11}^{\ell} \quad m_{211} \quad \rightarrow \quad \bar{m}'_{11}m_{11}^{\ell} \quad m_{11} \quad \rightarrow \quad \bar{m}'_{11} \quad \rightarrow \quad n'_{11}
 \end{array}$$

Notice that these determiners yield complete noun phrases which can be subject, object and attribute: *un bon vin blanc me plaît, j'aime un bon vin blanc, c'est un bon vin blanc.*

### Partitive complete noun phrase

French has complete noun phrases formed with the partitive article, *du, de la, de l', des*. Functioning as partitive <sup>6</sup>, *de* transforms an incomplete noun phrase into a complete one. This partitive noun phrase can be direct object of a verb (*Il mange du pain*), attribute (*C'est du sable*) or even subject (*des enfants jouent dans la rue*), i.e. the partitive article is understood as an indefinite article. In everyday French however, one rarely uses a noun phrase with the partitive article in the singular as the subject of a sentence: *\*Du pain est sur la table, ?De l'eau s'est infiltrée dans les fondements, ?Du sable gêne l'engrenage* are replaced by *il y a du pain sur la table, il y a de l'eau qui s'est infiltrée..., il y a du sable qui gêne l'engrenage.*

We introduce a new type  $\hat{n}_{gn}$ ,  $g = 1, 2; n = 1, 2$ , together with a super-type  $\hat{n}$  such that  $\hat{n}_{gn} \rightarrow \hat{n}$ . It is the type of a complete noun phrase which in general will not be used as subject. The plural partitive article *des*

<sup>6</sup> i.e. determining part of a mass or a group

transforms a plural count noun into a complete noun phrase. The same holds for the singular partitives *du*, *de la*, *de l'*, when preceding a mass-noun phrase.<sup>7</sup>

Hence the types

$$des : \hat{n}_{g2} \mathbf{x}_{g2}^{\ell}, \mathbf{x} = \mathbf{c}, \mathbf{c}'^8$$

$$du : \hat{n}_{11} \mathbf{m}_{11}^{\ell}$$

$$de : \hat{n}_{g1} \bar{\mathbf{m}}_{g1}^{\ell}$$

Examples

(Je mange) *des pommes*

$$\hat{n}_{22} \mathbf{c}_{22}^{\ell} \mathbf{c}_{22} \rightarrow \hat{n}_{22}$$

(Je mange) *du pain*

$$\hat{n}_{11} \mathbf{m}_{11}^{\ell} \mathbf{m}_{11} \rightarrow \hat{n}_{11}$$

(Il y a) *de l'air*

$$\hat{n}_{11} \bar{\mathbf{m}}_{11}^{\ell} \bar{\mathbf{m}}_{11} \rightarrow \hat{n}_{11}$$

Two comments: The first concerns the use of a partitive complete noun phrase as the subject of a sentence.

Compare

- (1) *Des gens vous demandent.*
- (2) \**Des nombres pairs sont divisibles par deux.*

The first sentence is generally accepted, see [Le bon usage], [Carrier], [Kleiber], whereas the second is rejected. The obvious reason lies in the meaning of the two sentences. The partitive article has the meaning of an existential quantifier, whereas the rejected sentence (2) would require a universal quantifier. The latter can be rendered by the definite article *les*:

- (3) *Les nombres pairs sont divisibles par deux.*

Our analysis assigns different types to the noun phrases *des nombres pairs* and *les nombres pairs*, namely  $\hat{n}_{12}$  and  $\mathbf{n}_{12}$ . By an appropriate typing of the French sentence structure, it will therefore be possible to accept (1) and (3) and to reject (2).

Our second comment concerns the use of *de* instead of *des* if the noun is preceded by an adjective. We have

- (i) *des fleurs*
- (ii) *des jolies fleurs*

<sup>7</sup> The partitive article is also used to construct a postnominal complement to a noun phrase like *le temps des cerises*, *le goût du pain*, *la forme de la pomme*. One would have to declare still more types for *des*, *du*, *de* which is beyond the scope of this paper.

<sup>8</sup> This implies that *des jolies fleurs* is accepted as a complete noun phrase.

(iii) *de jolies fleurs*

Our typing up to this point does not recognize (iii). This can easily be repaired by assigning to *de* new types, namely

$$de : \widehat{\mathbf{n}}_{g2} \mathbf{c}_{hg2}^\ell, h = 1, 2, \dots$$

$$d' : \widehat{\mathbf{n}}_{g2} \mathbf{c}'_{hg2}^\ell, h = 1, 2, \dots$$

Notice that this typing accepts *de jolies fleurs* and *d'autres fleurs* as complete noun phrases, but not *de fleurs*.

However, the latter will be analysed as a “quasi” complete noun phrase by assigning still more types to *de* as we shall explain next..

### Adverbial determiners

The preposition *de* preceded by adverbs of degree like *assez*, *beaucoup*, *combien*, *peu* also functions as a determiner, for example *peu d'enfants*, *beaucoup de pain*, *combien de sable*. Hence, *d'enfants*, *de pain*, *de sable* may be considered as “almost” complete noun phrases.

Yet another type  $\bar{\mathbf{x}}_{gn}$  will have to account for the “almost” complete noun phrases formed with *de* alone, without a subsequent definite article. Then we declare following new types for *de*

$$de : \bar{\mathbf{c}}_{gn} \mathbf{c}_{gn}^\ell, \bar{\mathbf{m}}_{g1} \mathbf{m}_{g1}^\ell,$$

$$d' : \bar{\mathbf{c}}_{gn} \mathbf{c}'_{gn}^\ell, \bar{\mathbf{m}}_{g1} \mathbf{m}'_{g1}^\ell,$$

$$assez, beaucoup : \mathbf{n}_{gn} \bar{\mathbf{x}}_{gn}^\ell, \text{ either } \mathbf{x} = \mathbf{c} \text{ and } n = 2 \text{ or } \mathbf{x} = \mathbf{m} \text{ and } n = 1.$$

Examples :

*(Je mange) beaucoup de pain*

$$\mathbf{n}_{11} \bar{\mathbf{m}}_{11}^\ell \quad \bar{\mathbf{m}}_{11} \mathbf{m}_{11}^\ell \quad \bar{\mathbf{m}}_{11} \quad \rightarrow \quad \mathbf{n}_{11}$$

*Beaucoup de chats (arrivent)*

$$\mathbf{n}_{12} \bar{\mathbf{c}}_{12}^\ell \quad \bar{\mathbf{c}}_{12} \mathbf{c}_{12}^\ell \quad \mathbf{c}_{12} \quad \rightarrow \quad \mathbf{n}_{12}$$

Anticipating on the analysis of the sentence structure, we have introduced two different types for the complete noun phrase,  $\mathbf{n}_{gn}$  and  $\widehat{\mathbf{n}}_{gn}$ . The subtypes  $\bar{\mathbf{x}}_{gn}$  of  $\mathbf{n}_{gn}$  were needed to recognise the transformation of a complete noun phrase formed with a definite article (*l'eau*) into a complete “indefinite” noun phrase (*de l'eau*).

## 4. Conclusion

In Section 1, we have established the conditions for extending a language fragment without changing the completeness and correctness of the typing of the original fragment. The results of Section 2 show how and when the efficient general type checking algorithm can be improved to a linear one; an important property in implementations. Sections 1 and 2 combine to an argument that complexity, besides the actual grammar, could and should influence the typing of language fragments. The criterion for linearity has already influenced our

typing in Section 3. There we have also explained how each position of an index in the subscript corresponds to a “feature”, used to organise the dictionary in a succinct way. Thus, meta-rules are not but part of the grammar, but help to organise the dictionary.

Future work should take advantage of this succinct organisation of the dictionary to devise efficient type assignment algorithms. Certainly, more criterions for language fragments which can be decided in linear time would be welcome. We are convinced that pregroups are the tool for real time applications verifying grammars. Larger and larger fragments of natural languages need to be typed.

## References

- [Barg Lamb] Danièle Bargelli, Joachim Lambek, An algebraic approach to the French sentence structure, Logical Aspects of Computational Linguistics, 4<sup>th</sup> International Conference, LACL 2001, Le Croisic, France, June 27-29, 2001, Proceedings 2099: pp. 62-78, 2001.
- [Carlier] Anne Carlier, La Résistance des articles *du* et *des* à l’interprétation générique, in D. Amiot et al., editors, Le syntagme nominal syntaxe et sémantique, Artois Presses Université, 2001
- [Earley] Jay Earley, An efficient context-free parsing algorithm, Communications of the AMC, Volume 13, Number 2, pp 94-102, 1970
- [Kleiber] Georges Kleiber, Indéfinis: lecture existentielle et lecture partitive, in G. Kleiber et al., editors, Typologie des groupes nominaux, Presses Universités Rennes, 2001
- [Lambek99] Joachim Lambek, Type Grammar revisited, in A. Lecomte et al., editors, Logical Aspects of Computational Linguistics, Springer LNAI 1582, pp.1 –27, 1999
- [Le bon usage] Maurice Grevisse, André Goosse, Le bon usage, grammaire française, Duculot, 2001