

# Parsing with Polymorphic Categorical Grammars

Matteo Capelletti<sup>1</sup> and Fabio Tamburini<sup>2</sup>

<sup>1</sup> Lix, École Polytechnique - France  
matteo.capelletti@elisanet.fi

<sup>2</sup> DSLO, University of Bologna - Italy  
fabio.tamburini@unibo.it

**Abstract.** In this paper we investigate the use of polymorphic categorical grammars as a model for parsing natural language. We will show that, despite the undecidability of the general model, a subclass of polymorphic categorical grammars, which we call *linear*, is mildly context-sensitive and we propose a polynomial parsing algorithm for these grammars.

## 1 Introduction

The simplest model of a categorical grammar is the so called Ajdukiewicz–Bar-Hillel calculus of [2] and [4]. Syntactic categories are formed from a given set of atoms as *functions*  $a/b$  and  $b\backslash a$ , with  $b$  and  $a$  categories. The intuitive meaning of a syntactic category of the form  $a/b$  (resp.  $b\backslash a$ ) is that it looks for an *argument* of category  $b$  to its right (resp. left) to give a category of type  $a$ . The resulting grammar system is known to be context-free.

Contemporary categorical grammars in the style of Ajdukiewicz–Bar-Hillel grammars are called *combinatory categorical grammars*, see [25]. Such systems adopt other forms of composition rules which enable them to generate non-context-free languages, see [29; 28]. The other main tradition of categorical grammar, the type-logical grammars of [20; 18], stemming from the work of [15], adopt special kinds of structural rules, that enable the system to generate non-context-free languages.

Both approaches increase the generative power of the basic system by adding special kinds of rules. In this paper, we adopt a different strategy which consists in keeping the elementary rule component of Ajdukiewicz–Bar-Hillel grammar and in introducing *polymorphic* categories, that is syntactic categories that contain variables ranging over categories. The inference process will be driven by unification, rather than by simple identity of formulas. We will see two kinds of polymorphic categorical grammars, one that is Turing complete and another, resulting from a restriction on the first, which is mildly context-sensitive. This second system, which is obviously the most interesting one for linguistics, has some important advantages with respect to the aforementioned ones. In respect to TLG, the polymorphic system we define is *polynomial*, as we will prove by providing a parsing algorithm. In respect to CCG (and most known TLG), our system is not affected by the so called *spurious ambiguity* problem, that is the problem of generating multiple, semantically equivalent, derivations.

A system somewhat similar to our polymorphic categorial grammar was studied in the late 80s as a kind of feature theoretic categorial syntax by [26] and [30], although, to our knowledge, after these works it has been almost completely neglected. While our system, in its most general form, allows the encoding of feature structures into categorial grammar, our primary concern will be to discuss its generative power and its computational properties. We will present general polymorphic categorial grammars, and show that they are undecidable. Then we will define a restricted class of such grammars which is mildly context-sensitive. We will show this by giving examples of the languages that can and cannot be generated, and by providing a polynomial parser for such grammars.

In relation to other mildly context-sensitive grammar formalism, the categorial grammar framework is well known for its close correspondence between syntactic derivation and the construction of semantic representation, see [27]. For this reason the extensions of the basic calculi with more powerful devices allows for the construction of sophisticated parsing systems dealing in parallel with both the syntax and the semantics of languages.

## 2 Polymorphic Ajdukiewicz–Bar-Hillel Grammar

In this paper we will assume that categories are not only of the form  $a/b$  or  $b \backslash a$ , but also  $a \otimes b$ , this is not standard in the categorial literature, but such product categories are useful as we will see in the examples that follow. From now on, we use capital letters  $A$ ,  $B$  and  $C$  and their variants as metavariables over formulas. As we said before, an expression of category  $A/B$  looks for an expression of category  $B$  to its right to give a compound expression of category  $A$  and dually for  $B \backslash A$ . Instead, an expression of category  $A \otimes B$  is an expression whose syntactic information is given by a component of category  $A$  and another of category  $B$ . We assume that the slashes bind more tightly than the product. We define a *sequent* as a pair where the first element, the antecedent, is a list of category, while the second, the succedent, contains a single category, and we write it as  $\Gamma \rightarrow A$ .

The deductive system given in Figure 1, which we call  $AB^\otimes$ , is a simple modification of the calculus of [14] to which it can easily be proved equivalent.

<b>Identity Axioms:</b>	$A \rightarrow A$
<b>Product Axioms:</b>	$A, B \rightarrow A \otimes B$
<b>Shifting Rules:</b>	$\frac{\Gamma \rightarrow C/A}{\Gamma, A \rightarrow C} (S_1) \qquad \frac{\Gamma \rightarrow A \backslash C}{A, \Gamma \rightarrow C} (S_2)$
<b>Cut Rules:</b>	$\frac{\Gamma, A \rightarrow C \quad \Delta \rightarrow A}{\Gamma, \Delta \rightarrow C} (C_1) \qquad \frac{\Gamma \rightarrow A \quad A, \Delta \rightarrow C}{\Gamma, \Delta \rightarrow C} (C_2)$

**Fig. 1.** Ajdukiewicz–Bar-Hillel calculus with product,  $AB^\otimes$ .

There are two ways in literature for extending the basic models of categorial grammars to generate non context-free languages. The first is the use of structural rules or other types of composition schemes. These approaches are characteristic of type-logical grammars, see [20; 18; 19], and combinatory categorial grammars (CCG), see [25; 3], and have been widely explored in the past. The second is based on the introduction of *polymorphism*. In this paper, we study this second approach.

The formalism of polymorphic categorial grammar that we are going to present is inspired by the polymorphic theory of types, see [17; 9]. Types may contain type variables together with constants, and these variables may be (implicitly or explicitly) quantified over. The idea of polymorphism is very simple and natural. Rather than defining a class of  $\text{id}$  functions  $\text{id}_{Int} :: Int \rightarrow Int$ ,  $\text{id}_{Char} :: Char \rightarrow Char$  and so forth, the function  $\text{id}$  is defined for any type  $\alpha$ , as  $\text{id} :: \forall \alpha. \alpha \rightarrow \alpha$  or  $\text{id} :: \alpha \rightarrow \alpha$  where  $\alpha$  is implicitly universally quantified.

The same idea is very natural also in linguistics, where, for example, coordination particles such as ‘and’ and ‘or’ are typically polymorphic, as they coordinate expressions of *almost* any syntactic category. Thus one can find in the categorial grammar literature several examples of polymorphic assignments for these expressions [15; 24; 8; 7].

Another example of Ajdukiewicz–Bar-Hillel style categorial grammars adopting a form of polymorphism are the unification categorial grammars of [26; 30; 10], where polymorphism is used at the level of feature structures.

In this section, we are going to explore Ajdukiewicz–Bar-Hillel categorial grammars *with product* ( $AB^\otimes$ , for short) extended with different forms of polymorphism. In the first place, we extend the notion of category as to include variables and define a sort of categorial grammar with ML style polymorphism, see [17; 12]. Variables are assumed to be implicitly universally quantified, with quantifiers in prenex position.

Consider the formula  $(\alpha \backslash \alpha) / \alpha$  which can be assigned to the word ‘and’ in a lexicon. Applied to ‘John’  $:: n$ , it will give the expression ‘and John’  $:: n \backslash n$ , while applied to ‘walks’  $:: n \backslash s$ , it will give the expression ‘and walks’  $:: (n \backslash s) \backslash (n \backslash s)$ , and so forth. Thus the right argument  $\alpha$  in  $(\alpha \backslash \alpha) / \alpha$  gets instantiated in the application process and the result of such instantiation, a *substitution*, is applied to the value  $\alpha \backslash \alpha$ . Hence the process of type inference for this kind of polymorphic categorial grammars requires nothing more than type unification and substitution. We call the resulting system Unification Ajdukiewicz–Bar-Hillel grammars with product,  $UAB^\otimes$  for short.

## 2.1 Unification Ajdukiewicz–Bar-Hillel Grammars

Syntactic categories of  $UAB^\otimes$  are defined as follows.

**Atoms:**  $\mathcal{A} ::= a, b, c, n, s, i \dots$   
**Variables:**  $\mathcal{V} ::= \alpha, \beta, \gamma \dots$   
**Categories:**  $\mathcal{F} ::= \mathcal{A} \mid \mathcal{V} \mid \mathcal{F} \otimes \mathcal{F} \mid \mathcal{F} \backslash \mathcal{F} \mid \mathcal{F} / \mathcal{F}$

Unification of categories is defined in (1). We use  $\sharp$  as a variable over  $\{/, \backslash, \otimes\}$  connectives. The *substitution* of a formula  $A$  for a variable  $\alpha$  in a formula  $B$ , denoted  $B[\alpha := A]$ , is defined as follows:

$$\begin{aligned} \alpha[\alpha := C] &= C & \alpha[\beta := C] &= \alpha \text{ if } \alpha \neq \beta \\ A[\alpha := C] &= A \text{ for } A \in \mathcal{A} & (A\sharp B)[\alpha := C] &= (A[\alpha := C]\sharp B[\alpha := C]). \end{aligned}$$

We write  $f \cdot g$  the composition of  $f$  and  $g$ , defined as  $\lambda x.f(g x)$ . Let  $V(B)$  be the set of variables occurring in  $B$ , we define the *unification* of  $A$  and  $B$ , denoted  $A \approx B$ , by the following recursion, which is taken with minor modifications from [5].

$$\begin{aligned} \alpha \approx B &= [\alpha := B] && \text{if } \alpha \notin V(B) \\ &= \text{Id} && \text{if } \alpha \equiv B \\ &= \text{fail} && \text{otherwise} \\ A\sharp B \approx A'\sharp B' &= (\sigma A \approx \sigma A') \cdot \sigma && \text{where } \sigma = B \approx B' \\ A \approx \alpha &= \alpha \approx A \end{aligned} \tag{1}$$

The unification Ajdukiewicz–Bar-Hillel calculus,  $\text{UAB}^\otimes$  is defined in Figure 2.<sup>3</sup>

<b>Identity Axioms:</b>	$A \rightarrow A$
<b>Product Axioms:</b>	$A, B \rightarrow A \otimes B$
<b>Shifting Rules:</b>	$\frac{\Gamma \rightarrow C/A}{\Gamma, A \rightarrow C} (S_1) \qquad \frac{\Gamma \rightarrow A \backslash C}{A, \Gamma \rightarrow C} (S_2)$
<b>Cut Rules:</b>	$\frac{\Gamma, A \rightarrow C \quad \Delta \rightarrow B}{\Gamma, \Delta \rightarrow C(A \approx B)} (C'_1) \qquad \frac{\Gamma \rightarrow B \quad A, \Delta \rightarrow C}{\Gamma, \Delta \rightarrow C(A \approx B)} (C'_2)$

**Fig. 2.** Unification Ajdukiewicz–Bar-Hillel calculus,  $\text{UAB}^\otimes$

We give here some examples of non context-free languages generated by  $\text{UAB}^\otimes$  grammars.

**Example 1** *We define the  $\text{UAB}^\otimes$  grammar for the language  $a^n b^n c^n$ ,  $n \geq 1$ , a well-known non-context-free language. Let grammar  $G_1$  consist of the following assignments:*

$$\begin{aligned} a &:: s/(b \otimes c) & b &:: b \\ a &:: (s/\alpha) \backslash (s/(b \otimes (\alpha \otimes c))) & c &:: c \end{aligned}$$

*We derive the string ‘aabbcc’. We write  $A$  for the formula  $(s/\alpha) \backslash (s/(b \otimes (\alpha \otimes c)))$ . For readability, boxes are drawn around the words that anchor the axioms to the lexicon.*

<sup>3</sup> Obviously, the rules involving unification are only defined if unification is defined.

$$\frac{\frac{\frac{\boxed{a}}{s/(b \otimes c)} \rightarrow s/(b \otimes c) \quad \frac{\boxed{a}}{A \rightarrow A} \quad \frac{\boxed{b}}{b \rightarrow b} \quad \frac{\boxed{c}}{c \rightarrow c}}{\frac{\boxed{a}}{s/\alpha, A \rightarrow s/(b \otimes (\alpha \otimes c))} \quad \frac{\boxed{b}}{b, c \rightarrow b \otimes c} \quad \frac{\boxed{c}}{c \rightarrow c}}{\frac{\boxed{a}}{s/(b \otimes c), A \rightarrow s/(b \otimes ((b \otimes c) \otimes c))} \quad \frac{\boxed{b}}{b, b, c, c \rightarrow b \otimes ((b \otimes c) \otimes c)}}{\frac{\boxed{a}}{s/(b \otimes c)}, A, b, b, c, c \rightarrow s}$$

To show that  $G_1$  indeed generates the language  $a^n b^n c^n$ ,  $n \geq 1$ , we proceed by induction on  $n$ . If  $n = 1$ , then ‘ $abc$ ’ is generated by axioms  $a/(b \otimes c) \rightarrow a/(b \otimes c)$ ,  $b \rightarrow b$  and  $c \rightarrow c$ . Assume that  $G_1$  generates  $a^n b^n c^n$ . Then  $G_1$  assigns  $a^n$  the category  $s/A$  for some  $A$  and  $b^n c^n$  the category  $A$ . We have  $a :: (s/\alpha) \setminus (s/(b \otimes (\alpha \otimes c)))$ . Hence,  $G_1$  assigns  $a^{n+1}$  the category  $s/(b \otimes (A \otimes c))$  and  $b^{n+1} c^{n+1}$  the category  $b \otimes (A \otimes c)$ . We conclude that  $G_1$  generates  $a^{n+1} b^{n+1} c^{n+1}$ .

**Example 2** We define a  $UAB^\otimes$  grammar for ‘ $ww$ ’,  $w \in \{a, b\}^+$ . Let grammar  $G_2$  consist of the following assignments:

$$\begin{array}{ll} a :: a & b :: b \\ a :: s/a & b :: s/b \\ a :: (s/\alpha) \setminus (s/(\alpha \otimes a)) & b :: (s/\alpha) \setminus (s/(\alpha \otimes b)) \end{array}$$

It is easy to see that grammar  $G_2$  generates exactly the language ‘ $ww$ ’ with  $w \in \{a, b\}^+$ . As in the case of  $G_1$ , type variables are used as accumulators for long-distance dependencies. Here we give an example deduction for the string ‘ $aabaab$ ’, using  $A$  as for  $(s/\alpha) \setminus (s/(\alpha \otimes a))$  and  $B$  for  $(s/\alpha) \setminus (s/(\alpha \otimes b))$ .

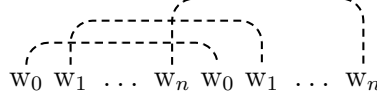
$$\frac{\frac{\frac{\boxed{a}}{s/a} \quad \frac{\boxed{a}}{s/\alpha, A \rightarrow s/(\alpha \otimes a)}}{\frac{\boxed{a}}{s/a, A \rightarrow s/((a \otimes a) \otimes b)}} \quad \frac{\frac{\boxed{b}}{B \rightarrow B}}{s/\alpha, B \rightarrow s/(\alpha \otimes b)} \quad \frac{\frac{\boxed{a}}{a \rightarrow a} \quad \frac{\boxed{a}}{a \rightarrow a} \quad \frac{\boxed{b}}{b \rightarrow b}}{\frac{\boxed{a}}{a, a \rightarrow (a \otimes a)} \quad \frac{\boxed{b}}{a, a, b \rightarrow (a \otimes a) \otimes b}}{\frac{\boxed{a}}{s/a, A, B \rightarrow s/((a \otimes a) \otimes b)} \quad \frac{\boxed{b}}{a, a, b \rightarrow (a \otimes a) \otimes b}}{\frac{\boxed{a}}{s/a, A, B, a, a, b \rightarrow s}}$$

A typical example of non context-freeness of natural language are the so called cross serial dependencies, which can be found, for instance, in Dutch subordinate clauses.

**Example 3** We define a  $UAB^\otimes$  grammar for Dutch cross-serial dependencies. An example is the following subordinate clause, from [25]:

*Ik Cecilia Henk de nijlpaarden zag helpen voeren.*  
 I Cecilia Henk the hippopotamuses saw help feed.  
 I saw Cecilia help Henk feed the hippopotamuses.

These constructs exhibit dependencies of the form ‘ $ww$ ’, where the  $i$ th words in the two halves are matched.



A sample lexicon generating the sentence in this example is the following:

<i>Ik, Cecilia, Henk, de nijlpaarden</i>	:: <i>n</i>
<i>zag</i>	:: $((n \otimes (n \otimes \alpha)) \setminus c) / (\alpha \setminus i)$
<i>helpen</i>	:: $((n \otimes \alpha) \setminus i) / (\alpha \setminus i)$
<i>voeren</i>	:: $n \setminus i$

With such a lexicon we obtain the following deduction for the subordinate clause (we write  $Z$  for the type of ‘zag’,  $H$  for that of ‘helpen’ and  $N$  for the string ‘Ik Cecilia Henk de nijlpaarden’).

$$\frac{\frac{\frac{\boxed{N}}{n \otimes (n \otimes (n \otimes n))} \quad \frac{\frac{\boxed{zag}}{Z \rightarrow Z} \quad \frac{Z, \alpha \setminus i \rightarrow (n \otimes (n \otimes \alpha)) \setminus c}{Z, (H, n \setminus i) \rightarrow (n \otimes (n \otimes (n \otimes n))) \setminus c}}{Z, (H, n \setminus i) \rightarrow (n \otimes (n \otimes (n \otimes n))) \setminus c} \quad \frac{\frac{\boxed{helpen}}{H \rightarrow H} \quad \frac{\boxed{voeren}}{n \setminus i \rightarrow n \setminus i}}{H, n \setminus i \rightarrow (n \otimes n) \setminus i}}{H, n \setminus i \rightarrow (n \otimes (n \otimes n)) \setminus c}}{n \otimes (n \otimes (n \otimes n)), (Z, (H, n \setminus i)) \rightarrow c}$$

These examples show that the languages generated by  $\text{UAB}^\otimes$  grammars properly include the context-free languages (since  $\text{AB}^\otimes$  grammars are instances of  $\text{UAB}^\otimes$  grammars).

With regard to the generative power of  $\text{UAB}^\otimes$  grammars, it can be proven that if we allow null assignments, that is assignments of the form  $\epsilon :: A$ , where  $\epsilon$  is the empty string, the  $\text{UAB}^\otimes$  formalism becomes undecidable. We can show this by translating in our categorial setting the argument of [13] to prove the Turing completeness of unification based attribute-value grammars. It is possible to encode a Turing machine  $M$  in a  $\text{UAB}^\otimes$  grammar  $G(M)$  such that  $G(M)$  generates the string ‘halt’ if and only if  $M$  halts with a blank input tape; this is enough to conclude the undecidability of  $\text{UAB}^\otimes$  grammars. Polymorphic null assignments, in fact, correspond quite neatly to lexical rules, as proposed in [6], leading to an undecidable formalism.

## 2.2 Constraining $\text{UAB}^\otimes$ Grammars

A constrain that we can impose on  $\text{UAB}^\otimes$  grammars to avoid undecidability is *linearity*. Roughly, we impose the restriction that any lexical type may contain at most one variable, occurring once in an argument position and once in value position. Thus,  $\alpha/\alpha$ ,  $(s/\alpha) \setminus (s/(\alpha \otimes a))$  are licit types, while  $(\alpha \setminus \alpha)/\alpha$ ,  $(s/(\alpha \otimes \beta)) \setminus (s/((\alpha \otimes \beta) \otimes a))$  and  $(s/(\alpha \otimes \alpha)) \setminus (s/((\alpha \otimes \alpha) \otimes a))$  are not. More precisely we define *linear* categories as the types  $F_2$  generated by the following context-free grammar.

$$\begin{array}{ll}
 \# ::= \otimes \mid / \mid \backslash & \#' ::= / \mid \backslash \\
 F_0 ::= A \mid F_0 \# F_0 & F_1 ::= F_1 \# F_0 \mid F_0 \# F_1 \mid \alpha \\
 F_2 ::= F_1 \# F_1 \mid F_0 \mid F_2 \# F_0 \mid F_0 \# F_2
 \end{array} \quad (2)$$

This definition of categories may deserve some comments. The interesting cases are the  $F_2$  formulas of the form  $A/B$  or  $B \backslash A$ , with  $A$  and  $B$  in  $F_1$  (the other cases are meant essentially to put these in context). Consider the case of  $A/B$ , then  $\alpha$  occurs exactly once in  $A$  and in  $B$ , since a  $F_1$  category contains the variable  $\alpha$  by construction. By analogy with lambda terms, we can think of the occurrence of  $\alpha$  in  $B$  as a *binder* (possibly a pattern-binder), and of the occurrence in  $A$  as the *bound* variable.

An  $\text{UAB}^\otimes$  grammar is linear if all its lexical assignments are linear. Furthermore, in linear  $\text{UAB}^\otimes$  grammar, we work by simple *variable instantiation*, rather than by a full-fledged unification algorithm. More precisely let us denote  $A^B$  a formula  $A$  with a distinguished occurrence of a subformula  $B$ .  $A^C$  is the formula obtained from  $A^B$  by replacing the occurrence of the subformula  $B$  with the formula  $C$ . The linear  $\text{UAB}^\otimes$  calculus results from the  $\text{UAB}^\otimes$  calculus in Figure 2 by replacing the Cut rules with the following *instantiation* rules.

$$\frac{\Delta \rightarrow A^B \quad \Gamma, A^\alpha \rightarrow C}{\Gamma, \Delta \rightarrow C[\alpha := B]} \quad \frac{\Gamma \rightarrow A^B \quad A^\alpha, \Delta \rightarrow C}{\Gamma, \Delta \rightarrow C[\alpha := B]} \quad (3)$$

Observe that given a linear  $\text{UAB}^\otimes$  grammar adopting the rules in 3, only linear types can occur in any of its deductions.

Observe also that the  $\text{UAB}^\otimes$  grammars for  $a^n b^n c^n$  and ww languages as well as that for the Dutch cross serial patterns, are all linear. On the other hand, no linear  $\text{UAB}^\otimes$  grammar can be given for the so called MIX or Bach language that is the language of the strings containing an equal number of a's, b's and c's<sup>4</sup>.

As we have the proper inclusion of context-free languages and the realization of limited cross-serial dependencies, in order to have a *mildly context-sensitive* grammar formalism we shall prove that linear  $\text{UAB}^\otimes$  grammars can be parsed in polynomial time. We do this in the next section by providing a parsing algorithm for linear  $\text{UAB}^\otimes$  grammars.

### 3 Polynomial Parsing with Linear $\text{UAB}^\otimes$ Grammars

In this section we define a polynomial parsing algorithm for linear  $\text{UAB}^\otimes$  grammars. It is based on a standard agenda-driven chart-parsing method and makes use of an external table, which we call *instantiation table*, for storing the 'partial' instantiations of variables. Let  $n$  be the length of the input string and  $Lex$

<sup>4</sup> To see this, we observe that the context-free language of the strings containing an equal number of a's and b's is not *linear*, in the sense of [11], see [16]. Hence for the MIX language, a  $\text{UAB}^\otimes$  grammar needs to bind two distinct variables for each symbol, what violates linearity.

the input lexicon. Cells of the instantiation table are denoted  $\mathcal{I}_{(i,k,j)}$ , where  $0 \leq i < j \leq n$  and  $0 \leq k \leq |Lex|$ . We extend the construction of formulas with two kinds of variables,  $\alpha_k$  and  $\alpha_{(i,k,j)}$  where  $i, k$  and  $j$  are as before. The difference between the two kinds of variables is that  $\alpha_k$  is an uninstantiated variable while  $\alpha_{(i,k,j)}$  is a variable  $\alpha_k$  which has been instantiated when an item spanning between  $i$  and  $j$  was generated. The algorithm assumes that different lexical entries contain different variables, that is for no  $k$  the variable  $\alpha_k$  occurs in two distinct lexical assignments<sup>5</sup>.

The parser operates on two kinds of items which we represent as

$$(i, \Delta \triangleright \Gamma \rightarrow A, j) \quad \text{and} \quad (i, \Gamma \triangleleft \Delta \rightarrow A, j)$$

where  $i$  and  $j$  are integers and  $\Delta \triangleright \Gamma \rightarrow A$  and  $\Gamma \triangleleft \Delta \rightarrow A$  are *sequents* in which exactly one occurrence of an auxiliary symbols (either  $\triangleright$  or  $\triangleleft$ ) appear. These symbols play a role similar to the dot in Earley parsing systems.

An item of the form  $(i, \Delta \triangleright \Gamma \rightarrow A, j)$  asserts that  $\Delta \Gamma \rightarrow A$  is derivable in linear  $\text{UAB}^{\otimes}$  and that  $w_{i+1} \dots w_j \Rightarrow^* \Delta$ . Furthermore, the items have a predictive component.

- In case  $A \equiv A' \otimes A''$  and  $\Delta \Gamma \equiv A' A''$ , it asserts that for some context  $A$ ,  $w_{l+1} \dots w_i A \Lambda \Rightarrow^* C$  with  $0 \leq l < i$ . This means that the item has been *predicted* from another item  $(l, \Xi \triangleright A \Lambda \rightarrow C, i)$ .
- In case  $A \equiv \alpha_{(i,k,j)}$  and  $\Delta \Gamma \equiv B$  for some formula  $B \in \mathcal{I}_{(i,k,j)}$ , it asserts that for some context  $A$ ,  $w_{l+1} \dots w_i A^{\alpha_{(i,k,j)}} \Lambda \Rightarrow^* C$  with  $0 \leq l < i$ . This means that the item has been *dereferenced* from another item  $(l, \Xi \triangleright A^{\alpha_{(i,k,j)}} \Lambda \rightarrow C, i)$ .

The dual conditions hold for  $(i, \Gamma \triangleleft \Delta \rightarrow A, j)$ . For simplicity, we write items of the form  $(i, \triangleleft \Delta \rightarrow A, j)$  and  $(i, \Delta \triangleright \rightarrow A, j)$  as  $(i, \Delta \rightarrow A, j)$ .

The proposed parsing algorithm for linear  $\text{UAB}^{\otimes}$  grammars is showed in Figure 3.

The correctness of the algorithm can be proven by adapting the methods of [1; 23] for the CYK and Earley parsers and by observing that the triple  $(i, k, j)$  resulting from the instantiation of a variable is determined by the unique name  $k$  of the variable  $\alpha_k$  and by the span  $(i, j)$  over which the instantiation has been determined. Observe that the Completion rules apply to a non-instantiated variable and therefore that instantiated variables behave like constants in the parsing process, whose value is determined by the Dereference rules.

To show that the resulting algorithm is polynomial, we follow the usual argument for agenda-driven chart-based parser evaluation, see for instance [22; 21]. The number of sequents that can occur in a cell of the chart for a given grammar with lexicon  $Lex$  is  $O(n^2 |Lex| |\Sigma|)$  where  $\Sigma$  is the set of subformulas of the lexicon<sup>6</sup> and  $|Lex|$  is the number of word-category pairs contained in the lexi-

<sup>5</sup> Clearly all these modifications do not affect linearity, and are motivated by correctness and efficiency reasons.

<sup>6</sup> That is for each lexical entry  $w :: A$ ,  $\Sigma$  contains all the subformulas of  $A$ . For example the formula  $(a \setminus b) / c$  generates the subformulas  $\{(a \setminus b) / c, a \setminus b, a, b, c\}$



*Input:* a string  $w = w_1 \dots w_n$  and an  $UAB^\otimes$  grammar  $G$ .

*Output:* Accept/reject.

*Data Structures:* an  $(n+1) \times (n+1)$  matrix  $\mathcal{T}$ , the *chart*, whose cells are sets of sequents, a set of items  $\mathcal{N}$ , the *agenda*, containing all the items to be processed and the instantiation table  $\mathcal{I}$  as defined before.

**Initialization:**

Let  $\mathcal{N} = \emptyset$  and  $\mathcal{T}_{(i,j)} = \emptyset \quad \forall i, j$ .

For  $i = 1$  to  $n$  do

$\mathcal{N} = \mathcal{N} \cup \{ (i-1, A \rightarrow A, i) \mid w_i :: A \in Lex \}$

**Main cycle:**

While  $\mathcal{N} \neq \emptyset$  do

remove one item  $\nu = (i, \Gamma \rightarrow A, j)$  from  $\mathcal{N}$ .

If  $\Gamma \rightarrow A \notin \mathcal{T}_{(i,j)}$ , then

add  $\Gamma \rightarrow A$  to chart  $\mathcal{T}_{(i,j)}$

**Shifting:**

If  $\nu = (i, \Gamma \rightarrow C/A, j)$ , then add  $(i, \Gamma \triangleright A \rightarrow C, j)$  to  $\mathcal{N}$ .

If  $\nu = (i, \Gamma \rightarrow A/C, j)$ , then add  $(i, A \triangleleft \Gamma \rightarrow C, j)$  to  $\mathcal{N}$ .

**Prediction:**

If  $\nu = (i, \Gamma \triangleright A \otimes B \Delta \rightarrow C, j)$ , then add  $(j, \triangleright A, B \rightarrow A \otimes B, j)$  to  $\mathcal{N}$ .

If  $\nu = (i, \Gamma A \otimes B \triangleleft \Delta \rightarrow C, j)$ , then add  $(i, A, B \triangleleft \rightarrow A \otimes B, i)$  to  $\mathcal{N}$ .

**$\epsilon$ -Scanning:**

If  $\nu = (i, \Gamma \triangleright A \Delta \rightarrow C, j)$  and  $\epsilon \Rightarrow^+ A$ , then add  $(i, \Gamma A \triangleright \Delta \rightarrow C, j)$  to  $\mathcal{N}$ .

If  $\nu = (i, \Gamma A \triangleleft \Delta \rightarrow C, j)$  and  $\epsilon \Rightarrow^+ A$ , then add  $(i, \Gamma \triangleleft A \Delta \rightarrow C, j)$  to  $\mathcal{N}$ .

**Completion:**

If  $\nu = (i, \Gamma \triangleright A \Delta \rightarrow C, j)$ , then for all  $\Lambda \rightarrow A \in \mathcal{T}_{(j,k)}$ , add  $(i, \Gamma A \triangleright \Delta \rightarrow C, k)$  to  $\mathcal{N}$ .

If  $\nu = (i, \Gamma A \triangleleft \Delta \rightarrow C, j)$ , then for all  $\Lambda \rightarrow A \in \mathcal{T}_{(k,i)}$ , add  $(k, \Gamma \triangleleft A \Delta \rightarrow C, j)$  to  $\mathcal{N}$ .

If  $\nu = (i, \Lambda \rightarrow A, j)$ , then  $\left\{ \begin{array}{l} \text{for all } \Gamma \triangleright A \Delta \rightarrow C \in \mathcal{T}_{(k,i)}, \text{ add } (k, \Gamma A \triangleright \Delta \rightarrow C, j) \text{ to } \mathcal{N}, \\ \text{for all } \Gamma A \triangleleft \Delta \rightarrow C \in \mathcal{T}_{(j,k)}, \text{ add } (i, \Gamma \triangleleft A \Delta \rightarrow C, k) \text{ to } \mathcal{N}. \end{array} \right.$

If  $\nu = (i, \Gamma \triangleright A^{\alpha_l} \Delta \rightarrow C, j)$ , then

for all  $\Lambda \rightarrow A^B \in \mathcal{T}_{(j,k)}$

add  $(i, \Gamma A^{\alpha_l} \triangleright \Delta \rightarrow C[\alpha_l := \alpha_{(i,l,k)}], k)$  to  $\mathcal{N}$  and update  $\mathcal{I}_{(i,l,k)} = \mathcal{I}_{(i,l,k)} \cup \{B\}$ .

If  $\nu = (i, \Gamma A^{\alpha_l} \triangleleft \Delta \rightarrow C, j)$ , then

for all  $\Lambda \rightarrow A^B \in \mathcal{T}_{(k,i)}$

add  $(k, \Gamma \triangleleft A^{\alpha_l} \Delta \rightarrow C[\alpha_l := \alpha_{(k,l,j)}], j)$  to  $\mathcal{N}$  and update  $\mathcal{I}_{(k,l,j)} = \mathcal{I}_{(k,l,j)} \cup \{B\}$ .

If  $\nu = (i, \Lambda \rightarrow A^B, j)$ , then

for all  $\Gamma \triangleright A^{\alpha_l} \Delta \rightarrow C \in \mathcal{T}_{(k,i)}$

add  $(k, \Gamma A^{\alpha_l} \triangleright \Delta \rightarrow C[\alpha_l := \alpha_{(k,l,j)}], j)$  to  $\mathcal{N}$  and update  $\mathcal{I}_{(k,l,j)} = \mathcal{I}_{(k,l,j)} \cup \{B\}$ .

for all  $\Gamma A^{\alpha_l} \triangleleft \Delta \rightarrow C \in \mathcal{T}_{(j,k)}$

add  $(i, \Gamma \triangleleft A^{\alpha_l} \Delta \rightarrow C[\alpha_l := \alpha_{(i,l,k)}], k)$  to  $\mathcal{N}$  and update  $\mathcal{I}_{(i,l,k)} = \mathcal{I}_{(i,l,k)} \cup \{B\}$ .

**Dereference:**

If  $\nu = (i, \Gamma \triangleright \alpha_{(k,l,m)} \Delta \rightarrow C, j)$  and  $A \in \mathcal{I}_{(k,l,m)}$  then add  $(j, \triangleright A \rightarrow \alpha_{(k,l,m)}, j)$  to  $\mathcal{N}$ .

If  $\nu = (i, \Gamma \alpha_{(k,l,m)} \triangleleft \Delta \rightarrow C, j)$  and  $A \in \mathcal{I}_{(k,l,m)}$  then add  $(j, A \triangleleft \rightarrow \alpha_{(k,l,m)}, j)$  to  $\mathcal{N}$ .

**Termination:** If, when  $\mathcal{N} = \emptyset$ ,  $\Gamma \rightarrow s \in \mathcal{T}_{(0,n)}$ , for some  $\Gamma$ , then accept, else reject.

**Fig. 3.** Recognition algorithm for linear  $UAB^\otimes$  grammars.

con. Then, the time is dominated by the Completion rules, that gives a global asymptotic complexity of  $O(n^5|Lex||\Sigma|)$ .

## 4 Conclusion

In this paper we have investigated some linguistic and computational properties of unification based categorial grammars. We have seen that, like other unification based grammar formalisms, unrestricted  $UAB^\otimes$  grammars are Turing complete. However, we have also seen that the constraint of *linearity* defines a linguistically interesting class of categorial grammars. Most notably, it locates the system among the mildly context-sensitive formalisms. Another pleasant aspect of the resulting system, at least with respect to any other CG-based mildly context-sensitive categorial formalism (be it CCG or type-logical grammar) is the absence of spurious ambiguity.

The work in [8] presents a decision procedure for a kind of polymorphic categorial grammar allowing at most two instances of the same variable in argument position and one in value (e.g.  $(X \setminus X)/X$ ). However, with such kind of polymorphism we can generate languages that are beyond the mildly context-sensitives (for instance, we can easily generate indexed languages such as  $\{www|w \in \{a,b\}^+\}$ ). This extension, and its implications on recognition complexity, are currently being investigated.

Despite the fact that in this paper we have been concerned only with a recognition algorithms, a proper *parsing* algorithm that provides all the parses of the input, handling also semantic information, can be easily provided. In that case, the correspondence between syntax and semantics typical of categorial grammars and the absence of spurious ambiguity of our systems become important ingredients of the resulting natural language parser.

We wish to conclude this section by observing that the linearity constraint can also be relaxed. For instance, while preserving the condition that only one variable occurs in a formula, we can allow more than two occurrences of this variable. Then bounded languages such as  $w^i$  or  $a_1^i a_2^i \dots a_n^i$ , which are generalizations of ‘ww’ and  $a^i b^i c^i$ , can easily be generated and recognized in polynomial time.

## References

1. A. Aho and J. Ullman. *The Theory of Parsing, Translation and Compiling*, volume 1: Parsing. Prentice-Hall, INC., 1972.
2. K. Ajdukiewicz. Die syntaktische Konnexität. *Studia Philosophica*, 1:1–27, 1935.
3. J. Baldridge. *Lexically Specified Derivational Control in Combinatory Categorical Grammar*. PhD thesis, University of Edinburgh, 2002.
4. Y. Bar-Hillel. A quasi-arithmetical notation for syntactic description. *Language*, 29:47–58, 1953.
5. H. Barendregt. Lambda calculus with types. In S. Abramsky, Dov M. Gabbay, and T. S. E. Maibaum, editors, *Handbook of Logic in Computer Science*, volume 2, pages 117–309. Oxford University Press, 1992.
6. B. Carpenter. The generative power of categorial grammars and head-driven phrase structure grammars with lexical rules. *Computational Linguistics*, 17(3):301–313, 1991.
7. S. Clark and J. R. Curran. Wide-coverage efficient statistical parsing with ccg and log-linear models. *Computational Linguistics*, 33(4):493–552, 2007.
8. M. Emms. Parsing with polymorphism. In *EACL*, pages 120–129, 1993.
9. J.-Y. Girard, P. Taylor, and Y. Lafont. *Proofs and Types*. Cambridge University Press, 1989.
10. D. Heylen. *Types and Sorts. Resource logic for feature checking*. PhD thesis, UiL-OTS, Utrecht, 1999.
11. J. E. Hopcroft and J. D. Ullman. *Introduction to Automata Theory, Languages and Computation*. Addison-Wesley, 1979.
12. G. Huet. A uniform approach to type theory. In *Logical foundations of functional programming*, pages 337–397. Addison-Wesley Longman Publishing Co., Inc., Boston, MA, USA, 1990.
13. M. Johnson. *Attribute-Value Logic and the Theory of Grammar*, volume 16 of *CSLI Lecture Notes*. CSLI, Stanford, California, 1988.
14. M. Kandulski. The equivalence of nonassociative Lambek categorial grammars and context-free grammars. *Zeitschrift für Mathematische Logik und Grundlagen der Mathematik*, 34:41–52, 1988.
15. J. Lambek. The mathematic of sentence structure. *American Mathematical Monthly*, 65(3):154–170, 1958.
16. P. Linz. *An introduction to formal languages and automata*. D. C. Heath and Company, Lexington, MA, USA, 1990.
17. R. Milner. A theory of type polymorphism in programming. *Journal of Computer and System Sciences*, 17:348–375, 1978.
18. M. Moortgat. Categorical type logics. In J. van Benthem and A. ter Meulen, editors, *Handbook of Logic and Language*, pages 93–177. Elsevier, Amsterdam, 1997.
19. R. Moot. *Proof Nets for Linguistic Analysis*. PhD thesis, UiL-OTS, Utrecht, 2002.

20. G. Morrill. *Type Logical Grammar: Categorical Logic of Signs*. Kluwer, Dordrecht, 1994.
21. M.-J. Nederhof and G. Satta. Tabular parsing. In C. Martin-Vide, V. Mitran, and G. Paun, editors, *Formal Languages and Applications, Studies in Fuzziness and Soft Computing 148*, pages 529–549. Springer, 2004.
22. G. Satta and O. Stock. Bidirectional context-free grammar parsing for natural language processing. *AIJ: Artificial Intelligence*, 69, 1994.
23. K. Sikkil. Parsing schemata and correctness of parsing algorithms. *Theoretical Computer Science*, 199, 1998.
24. M. Steedman. Dependency and coordination in the grammar of dutch and english. *Language*, 61(3):523–568, 1985.
25. M. Steedman. *The Syntactic Process*. The MIT Press, 2000.
26. Hans Uszkoreit. Categorical unification grammars. In *COLING*, pages 187–194, 1986.
27. J. van Benthem. *Essays in Logical Semantics*. Reidel, Dordrecht, 1986.
28. K. Vijay-Shanker and D. J. Weir. The equivalence of four extensions of context-free grammars. *Mathematical Systems Theory*, 27, 1994.
29. D. J. Weir and A. K. Joshi. Combinatory categorical grammars: Generative power and relationship to linear context-free rewriting systems. In *Meeting of the Association for Computational Linguistics*, pages 278–285, 1988.
30. H. Zeevat. Combining categorical grammar and unification. In U. Reyle and C. Rohrer, editors, *Natural Language Parsing and Linguistic Theories*, pages 202–229. D. Reidel, Dordrecht, 1988.