

# Dependency Constraints for Lexical Disambiguation

**Guillaume Bonfante**

LORIA INPL

guillaume.bonfante@loria.fr

**Bruno Guillaume**

LORIA INRIA

bruno.guillaume@loria.fr

**Mathieu Morey**

LORIA Nancy-Université

mathieu.morey@loria.fr

## Abstract

We propose a generic method to perform lexical disambiguation in lexicalized grammatical formalisms. It relies on dependency constraints between words. The soundness of the method is due to invariant properties of the parsing in a given grammar that can be computed statically from the grammar.

## 1 Introduction

In this work, we propose a method of lexical disambiguation based on the notion of dependencies.

In modern linguistics, Lucien Tesnière developed a formal and sophisticated theory with dependencies (Tesnière, 1959). Nowadays, many current grammatical formalisms rely more or less explicitly on the notion of dependencies between words. The most straightforward examples are formalisms in the Dependency Grammars family but it is also true of the phrase structure based formalisms which consider that words introduce incomplete syntactic structures which must be completed by other words. This idea is at the core of Categorical Grammars (CG) (Lambek, 1958) and all its trends such as Abstract Categorical Grammars (ACG) (de Groote, 2001) or Combinatory Categorical Grammars (CCG) (Steedman, 2000), being mostly encoded in their type system. Dependencies in CG were studied in (Moortgat and Morrill, 1991) and for CCG in (Clark et al., 2002; Koller and Kuhlmann, 2009). Other formalisms can be viewed as modeling and using dependencies, such as Tree Adjoining Grammars (TAG) (Joshi, 1987) with their substitution and adjunction operations. Dependencies for TAG were studied in (Joshi and Rambow, 2003). More recently, Marchand et al. (2009) showed that it is also possible to extract a dependency structure from a syntactic analysis in Interaction Grammars (IG) (Guillaume and Perrier, 2008).

Another much more recent concept of polarity can be used in grammatical formalisms to express that words introduce incomplete syntactic structures. IG directly use polarities to describe these structures but it is also possible to use polarities in other formalisms in order to make explicit the more or less implicit notion of incomplete structures: for instance, in CG (Lamarche, 2008) or in TAG (Kahane, 2006; Bonfante et al., 2004; Gardent and Kow, 2005). On this regard, Marchand et al. (2009) exhibited a direct link between polarities and dependencies. This encourages us to say that in many respects dependencies and polarities are two sides of the same coin.

The aim of this paper is to show that dependencies can be used to express constraints on the taggings of a sentence and hence these dependency constraints can be used to partially disambiguate the words of a sentence. We will see that, in practice, using the link between dependencies and polarities, these dependency constraints can be computed directly from polarized structures.

Exploiting the dependencies encoded in lexical entries to perform disambiguation is the intuition behind supertagging (Bangalore and Joshi, 1999), a method introduced for LTAG and successfully applied since then to CCG (Clark and Curran, 2004) and HPSG (Ninomiya et al., 2006). These approaches select the most likely lexical entry (entries) for each word, based on Hidden Markov Models or Maximum Entropy Models. Like the work done by Boullier (2003), our method is not based on statistics nor heuristics, but on a necessary condition of the deep parsing. Consequently, we accept to have more than one lexical tagging for a sentence, as long as we can ensure to have the good ones (when they exist!). This property is particularly useful to ensure that the deep parsing will not fail because of an error at the disambiguation step.

In wide-coverage lexicalized grammars, a word

typically has about 10 corresponding lexical descriptions, which implies that for a short sentence of 10 words, we get  $10^{10}$  possible taggings. It is not reasonable to treat them individually. To avoid this, it is convenient to use an automaton to represent the set of all paths. This automaton has linear size with regard to the initial lexical ambiguity. The idea of using automata is not new. In particular, methods based on Hidden Markov Models (HMM) use such a technique for part-of-speech tagging (Kupiec, 1992; Merialdo, 1994). Using automata, we benefit from dynamic programming procedures, and consequently from an exponential temporal and space speed up.

## 2 Abstract Grammatical Framework

Our filtering method is applicable to any lexicalized grammatical formalism which exhibits some basic properties. In this section we establish these properties and define from them the notion of Abstract Grammatical Framework (AGF).

Formally, an **Abstract Grammatical Framework** is an  $n$ -tuple  $(\mathcal{V}, \mathcal{S}, \mathcal{G}, \mathbf{anc}, \mathcal{F}, \mathbf{p}, \mathbf{dep})$  where:

- $\mathcal{V}$  is the **vocabulary**: a finite set of words of the modeled natural language;
- $\mathcal{S}$  is the set of **syntactic structures** used by the formalism;
- $\mathcal{G} \subset \mathcal{S}$  is the **grammar**: the finite set of initial syntactic structures; a finite list  $[t_1, \dots, t_n]$  of elements of  $\mathcal{G}$  is called a **lexical tagging**;
- $\mathbf{anc} : \mathcal{G} \rightarrow \mathcal{V}$  maps initial syntactic structures to their anchors;
- $\mathcal{F} \subset \mathcal{S}$  is the set of final syntactic structures that the parsing process builds (for instance trees);
- $\mathbf{p}$  is the **parsing function** from lexical taggings to finite subsets of  $\mathcal{F}$ ;
- $\mathbf{dep}$  is the **dependency function** which maps a couple composed of a lexical tagging and a final syntactic structure to dependency structures.

Note that the  $\mathbf{anc}$  function implies that the grammar is lexicalized: each initial structure in  $\mathcal{G}$  is associated to an element of  $\mathcal{V}$ . Note also that no particular property is required on the dependency

structures that are obtained with the  $\mathbf{dep}$  function, they can be non-projective for instance.

We call **lexicon** the function (written  $\ell$ ) from  $\mathcal{V}$  to subsets of  $\mathcal{G}$  defined by:

$$\ell(w) = \{t \in \mathcal{G} \mid \mathbf{anc}(t) = w\}.$$

We will say that a lexical tagging  $L = [t_1, \dots, t_n]$  is a lexical tagging of the sentence  $[\mathbf{anc}(t_1), \dots, \mathbf{anc}(t_n)]$ .

The final structures in  $\mathbf{p}(L) \subset \mathcal{F}$  are called the **parsing solutions** of  $L$ .

Henceforth, in our examples, we will consider the ambiguous French sentence (1).

(1) “*La belle ferme la porte*”

**Example 1** We consider the following toy AGF, suited for parsing our sentence:

- $\mathcal{V} = \{“la”, “belle”, “ferme”, “porte”\}$ ;
- the grammar  $\mathcal{G}$  is given in the table below: each  $\times$  corresponds to an element in  $\mathcal{G}$ , written with the category and the French word as subscript. For instance, the French word “porte” can be either a common noun (“door”) or a transitive verb (“hangs”); hence  $\mathcal{G}$  contains the 2 elements  $\text{CN}_{\text{porte}}$  and  $\text{TrV}_{\text{porte}}$ .

	<i>la</i>	<i>belle</i>	<i>ferme</i>	<i>porte</i>
Det	×			
LAdj		×	×	
RAAdj		×	×	
CN	×	×	×	×
Clit	×			
TrV			×	×
IntrV			×	

In our example, categories stand for, respectively: determiner, left adjective, right adjective, common noun, clitic pronoun, transitive verb and intransitive verb.

With respect to our lexicon, for sentence (1), there are  $3 \times 3 \times 5 \times 3 \times 2 = 270$  lexical taggings.

The parsing function  $\mathbf{p}$  is such that 3 lexical taggings have one solution and the 267 remaining ones have no solution; we do not need to precise the final structures, so we only give the English translation as the result of the parsing function:

- $\mathbf{p}([\text{Det}_{la}, \text{CN}_{belle}, \text{TrV}_{ferme}, \text{Det}_{la}, \text{CN}_{porte}]) = \{\text{“The nice girl closes the door”}\}$
- $\mathbf{p}([\text{Det}_{la}, \text{LAdj}_{belle}, \text{CN}_{ferme}, \text{Clit}_{la}, \text{TrV}_{porte}]) = \{\text{“The nice farm hangs it”}\}$
- $\mathbf{p}([\text{Det}_{la}, \text{CN}_{belle}, \text{RAdj}_{ferme}, \text{Clit}_{la}, \text{TrV}_{porte}]) = \{\text{“The firm nice girl hangs it”}\}$

### 3 The Companionship Principle

We have stated in the previous section the framework and the definitions required to describe our principle.

#### 3.1 Potential Companion

We say that  $u \in \mathcal{G}$  is a **companion** of  $t \in \mathcal{G}$  if  $\text{anc}(t)$  and  $\text{anc}(u)$  are linked by a dependency in  $\text{dep}(L, m)$  for some lexical tagging  $L$  which contains  $t$  and  $u$  and some  $m \in \mathbf{p}(L)$ . The subset of elements of  $\mathcal{G}$  that are companions of  $t$  is called the **potential companion set** of  $t$ .

The Companionship Principle says that if a lexical tagging contains some  $t$  but no potential companion of  $t$ , then it can be removed.

In what follows, we will generalize a bit this idea in two ways. First, the same  $t$  can be implied in more than one kind of dependency and hence it can have several different companion sets with respect to the different kinds of dependencies. Secondly, it can be the case that some companion  $t$  has to be on the right (resp. on the left) to fulfill its duty. These generalizations are done through the notion of atomic constraints defined below.

#### 3.2 Atomic constraints

We say that a pair  $(\mathcal{L}, \mathcal{R})$  of subsets of  $\mathcal{G}$  is an **atomic constraint** for an initial structure  $t \in \mathcal{G}$  if for each lexical tagging  $L = [t_1, \dots, t_n]$  such that  $\mathbf{p}(L) \neq \emptyset$  and  $t = t_i$  for some  $i$  then:

- either there is some  $j < i$  such that  $t_j \in \mathcal{L}$ ,
- or there is some  $j > i$  such that  $t_j \in \mathcal{R}$ .

In other words,  $(\mathcal{L}, \mathcal{R})$  lists the potential companions of  $t$ , respectively on the left and on the right.

A **system of constraints** for a grammar  $\mathcal{G}$  is a function  $\mathcal{C}$  which associates a finite set of atomic constraints to each element of  $\mathcal{G}$ .

The Companionship Principle is an immediate consequence of the definition of atomic constraints. It can be stated as the necessary condition:

#### The Companionship Principle

If a lexical tagging  $[t_1, \dots, t_n]$  has a solution then for all  $i$  and for all atomic constraints  $(\mathcal{L}, \mathcal{R}) \in \mathcal{C}(t_i)$

- $\{t_1, \dots, t_{i-1}\} \cap \mathcal{L} \neq \emptyset$
- or  $\{t_{i+1}, \dots, t_n\} \cap \mathcal{R} \neq \emptyset$ .

**Example 2** Often, constraints can be expressed independently of the anchors. In our example, we use the category to refer to the subset of  $\mathcal{G}$  of structures defined with this category:  $\text{LAdj}$  for instance refers to the subset  $\{\text{LAdj}_{belle}, \text{LAdj}_{ferme}\}$ .

We complete the example of the previous section with the following constraints<sup>1</sup>:

❶	$t \in \text{CN} \Rightarrow (\text{Det}, \emptyset) \in \mathcal{C}(t)$
❷	$t \in \text{LAdj} \Rightarrow (\emptyset, \text{CN}) \in \mathcal{C}(t)$
❸	$t \in \text{RAdj} \Rightarrow (\text{CN}, \emptyset) \in \mathcal{C}(t)$
❹	$t \in \text{Det} \Rightarrow (\emptyset, \text{CN}) \in \mathcal{C}(t)$
❺	$t \in \text{Det} \Rightarrow (\text{TrV}, \text{TrV} \cup \text{IntrV}) \in \mathcal{C}(t)$
❻	$t \in \text{TrV} \Rightarrow (\text{Clit}, \text{Det}) \in \mathcal{C}(t)$
❼	$t \in \text{TrV} \Rightarrow (\text{Det}, \emptyset) \in \mathcal{C}(t)$
❽	$t \in \text{IntrV} \Rightarrow (\text{Det}, \emptyset) \in \mathcal{C}(t)$
❾	$t \in \text{Clit} \Rightarrow (\emptyset, \text{TrV}) \in \mathcal{C}(t)$

The two constraints ❹ and ❺ for instance express that every determiner is implied in two dependencies. First, it must find a common noun on its right to build a noun phrase. Second, the noun phrase has to be used in a verbal construction.

Now, let us consider the lexical tagging:  $[\text{Det}_{la}, \text{LAdj}_{belle}, \text{TrV}_{ferme}, \text{Clit}_{la}, \text{CN}_{porte}]$  and the constraint ❾ (a clitic is waiting for a transitive verb on its right). This constraint is not fulfilled by the tagging so this tagging has no solution.

#### 3.3 The “Companionship Principle” language

Actually, a lexical tagging is an element of the formal language  $\mathcal{G}^*$  and we can consider the following three languages. First,  $\mathcal{G}^*$  itself. Second, the set  $C \subseteq \mathcal{G}^*$  corresponds to the lexical taggings which can be parsed. The aim of lexical disambiguation is then to exhibit for each sentence  $[w_1, \dots, w_n]$  all the lexical taggings that are within  $C$ . Third, the Companionship Principle defines the language  $P$  of lexical taggings which verify this Principle.  $P$  squeezes between the two lat-

<sup>1</sup>These constraints are relative to our toy grammar and are not linguistically valid in a larger context.

ter sets  $C \subseteq P \subseteq \mathcal{G}^*$ . Remarkably, the language  $P$  can be described as a regular language. Since  $C$  is presumably not a regular language (at least for natural languages!),  $P$  is a better regular approximation than the trivial  $\mathcal{G}^*$ .

Let us consider one lexical entry  $t$  and an atomic constraint  $(\mathcal{L}, \mathcal{R}) \in \mathcal{C}(t)$ . Then, the set of lexical taggings verifying this constraint can be described as

$$L_{t:(\mathcal{L}, \mathcal{R})} = \mathcal{C}((\mathcal{C}\mathcal{L})^*t(\mathcal{C}\mathcal{R})^*)$$

where  $\mathcal{C}$  denotes the complement of a set.

Since  $P$  is defined as the lexical taggings verifying all constraints,  $P$  is a regular language defined by :

$$P = \bigcap_{(\mathcal{L}, \mathcal{R}) \in \mathcal{C}(t)} L_{t:(\mathcal{L}, \mathcal{R})}$$

From the Companionship Principle, we derive a lexical disambiguation Principle which simply tests tagging candidates with  $P$ . Notice that  $P$  can be statically computed (at least, in theory) from the grammar itself.

**Example 3** For instance, for our example grammar, this automaton is given in the figure 1 where  $c=Clit$ ,  $n=CN$ ,  $d=Det$ ,  $i=IntrV$ ,  $l=LAdj$ ,  $r=RAdj$  and  $t=TrV$ .

A rough approximation of the size of the automaton corresponding to  $P$  can be easily computed. Since each automaton  $L_{t:(\mathcal{L}, \mathcal{R})}$  has 4 states,  $P$  has at most  $4^m$  states where  $m$  is the number of atomic constraints. For instance, the grammar used in the experiments contains more than one atomic constraint for each lexical entry, and  $m > |\mathcal{G}| > 10^6$ . Computing  $P$  by brute-force is then intractable.

## 4 Implementation of the Companionship Principle with automata

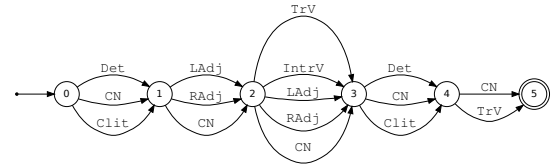
In this section we show how to use the Companionship Principle for disambiguation. Actually, we propose two implementations based on the principle, an exact one and an approximate one. The latter is really fast and can be used as a first step before applying the first one.

### 4.1 Automaton to represent sets of lexical taggings

The number of lexical taggings grows exponentially with the length of sentences. To avoid that, we represent sets of lexical taggings as the sets of paths of some acyclic automata where transitions

are labeled by elements of  $\mathcal{G}$ . We call such an automaton a **lexical taggings automaton (LTA)**. Generally speaking, such automata save a lot of space. For instance, given a sentence  $[w_1, \dots, w_n]$  the number of lexical taggings to consider at the beginning of the parsing process is  $\prod_{1 \leq i \leq n} |\ell(w_i)|$ . This set of taggings can be efficiently represented as the set of paths of the automaton with  $n + 1$  states  $s_0, \dots, s_n$  and with a transition from  $s_{i-1}$  to  $s_i$  with the label  $t$  for each  $t \in \ell(w_i)$ . This automaton has  $\sum_{1 \leq i \leq n} |\ell(w_i)|$  transitions.

**Example 4** With the data of the previous examples, we have the initial automaton:



To improve readability, only the categories are given on the edges, while the French words can be inferred from the position in the automaton.

### 4.2 Exact Companionship Principle (ECP)

Suppose we have a LTA  $\mathcal{A}$  for a sentence  $[w_1, \dots, w_n]$ . For each transition  $t$  and for each atomic constraint in  $(\mathcal{L}, \mathcal{R}) \in \mathcal{C}(t)$ , we construct an automaton  $\mathcal{A}_{t, \mathcal{L}, \mathcal{R}}$  in the following way.

Each state  $s$  of  $\mathcal{A}_{t, \mathcal{L}, \mathcal{R}}$  is labeled with a triple composed of a state of the automaton  $\mathcal{A}$  and two booleans. The intended meaning of the first boolean is to say that each path reaching this state passes through the transition  $t$ . The second boolean means that the atomic constraint  $(\mathcal{L}, \mathcal{R})$  is necessarily fulfilled.

The initial state is labeled  $(s_0, F, F)$  where  $s_0$  is the initial state of  $\mathcal{A}$  and other states are labeled as follows: if  $s \xrightarrow{u} s'$  in  $\mathcal{A}$  then, in  $\mathcal{A}_{t, \mathcal{L}, \mathcal{R}}$ , we have:

1.  $(s, F, b) \xrightarrow{u} (s', T, b)$  if  $u = t$
2.  $(s, F, b) \xrightarrow{u} (s', F, T)$  if  $u \in \mathcal{L}$
3.  $(s, F, b) \xrightarrow{u} (s', F, b)$  if  $u \notin \mathcal{L}$
4.  $(s, T, b) \xrightarrow{u} (s', T, T)$  if  $u \in \mathcal{R}$
5.  $(s, T, b) \xrightarrow{u} (s', T, b)$  if  $u \notin \mathcal{R}$

where  $b \in \{T, F\}$ .



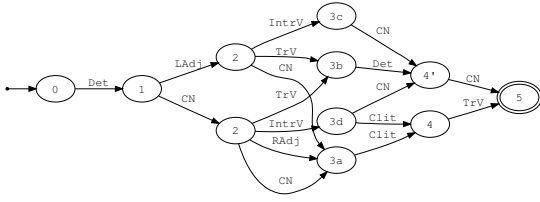
That is, we output :

$$\mathcal{A}_{CP} = \bigcap_{1 \leq i \leq n, t \in \mathcal{A}; (\mathcal{L}, \mathcal{R}) \in \mathcal{C}(t)} \mathcal{A}_{t, \mathcal{L}, \mathcal{R}}$$

It can be shown that the automaton is the same as the one obtained by intersection with the automaton of the language defined in 3.3:

$$\mathcal{A}_{CP} = \mathcal{A} \cap P$$

**Example 6** In our example, the intersection of the 9 automata built for the atomic constraints is given below:



This automaton has 8 paths: there are 8 lexical taggings which fulfill every constraint.

### 4.3 Approximation: the Quick Companionship Principle (QCP)

The issue with the previous algorithm is that it involves a large number of automata (actually  $O(n)$ ) where  $n$  is the size of the input sentence. Each of these automata has size  $O(n)$ . The theoretical complexity of the intersection is then  $O(n^n)$ . Sometimes, we face the exponential. So, let us provide an algorithm which approximates the Principle. The idea is to consider at the same time all the paths that contain some transition.

We consider a LTA  $\mathcal{A}$ . We write  $\prec_{\mathcal{A}}$  the precedence relation on transitions in an automaton  $\mathcal{A}$ . We define  $l_{\mathcal{A}}(t) = \{u \in \mathcal{G}, u \prec_{\mathcal{A}} t\}$  and  $r_{\mathcal{A}}(t) = \{u \in \mathcal{G}, t \prec_{\mathcal{A}} u\}$ .

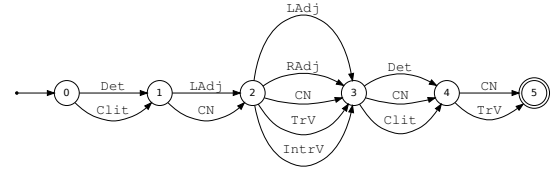
For each transition  $s \xrightarrow{t} s'$  and each constraint  $(\mathcal{L}, \mathcal{R}) \in \mathcal{C}(t)$ , if  $l_{\mathcal{A}}(t) \cap \mathcal{L} = \emptyset$  and  $r_{\mathcal{A}}(t) \cap \mathcal{R} = \emptyset$ , then none of the lexical taggings which use the transition  $t$  has a solution and the transition  $t$  can be safely removed from the automaton.

This can be computed by a double-for loop: for each atomic constraint of each transition, verify that either the left context or the right context of the transition contains some structure to solve the constraint. Observe that the cost of this algorithm is  $O(n^2)$ , where  $n$  is the size of the input automaton.

Note that one must iterate this algorithm until a fixpoint is reached. Indeed, removing a transition

which serves as a potential companion breaks the verification. Nevertheless, since for each step before the fixpoint is reached, we remove at least one transition, we iterate the double-for at most  $O(n)$  times. The complexity of the whole algorithm is then  $O(n^3)$ . In practice, we have observed that the complexity is closer to  $O(n^2)$ : only 2 or 3 loops are enough to reach the fixpoint.

**Example 7** If we apply the QCP to the automaton of Example 4, in the first step, only the transition  $0 \xrightarrow{\text{CN}} 1$  is removed by applying the atomic constraint ❶. In the next step, the transition  $1 \xrightarrow{\text{RAdj}} 2$  is removed by applying the atomic constraint ❸. The fixpoint is reached and the output automaton (with 120 paths) is:



## 5 The Generalized Companionship Principle

In practice, of course, we have to face the problem of the computation of the constraints. In a large coverage grammar, the size of  $\mathcal{G}$  is too big to compute all the constraints in advance. However, as we have seen in example 2 we can identify subsets of  $\mathcal{G}$  that have the same constraints; the same way, we can use these subsets to give a more concise presentation of the  $\mathcal{L}$  and  $\mathcal{R}$  sets of the atomic constraints. This motivates us to define a Generalized Principle which is stated on a quotient set of  $\mathcal{G}$ .

### 5.1 Generalized atomic constraints

Let  $\mathcal{U}$  be a set of subsets of  $\mathcal{G}$  that are a partition of  $\mathcal{G}$ . For  $t \in \mathcal{G}$ , we write  $\bar{t}$  the subset of  $\mathcal{U}$  which contains  $t$ .

We say that a pair  $(\mathcal{L}, \mathcal{R})$  of subsets of  $\mathcal{U}$  is a **generalized atomic constraint** for  $u \in \mathcal{U}$  if for each lexical tagging  $L = [t_1, \dots, t_n]$  such that  $\mathbf{p}(L) \neq \emptyset$  and  $u = \bar{t}_i$  for some  $i$  then:

- either there is some  $j < i$  such that  $\bar{t}_j \in \mathcal{L}$ ,
- or there is some  $j > i$  such that  $\bar{t}_j \in \mathcal{R}$ .

A **system of generalized constraints** for a partition  $\mathcal{U}$  of a grammar  $\mathcal{G}$  is a function  $\mathcal{C}$  which asso-

ciates a finite set of generalized atomic constraints to each element of  $\mathcal{U}$ .

## 5.2 The Generalized Principle

The Generalized Companionship Principle is then an immediate consequence of the previous definition and can be stated as the necessary condition:

### The Generalized Companionship Principle

If a lexical tagging  $[t_1, \dots, t_n]$  has a solution then for all  $i$  and for all generalized atomic constraints  $(\mathcal{L}, \mathcal{R}) \in \mathcal{C}(t_i)$

- $\{\overline{t_1}, \dots, \overline{t_{i-1}}\} \cap \mathcal{L} \neq \emptyset$
- or  $\{\overline{t_{i+1}}, \dots, \overline{t_n}\} \cap \mathcal{R} \neq \emptyset$ .

**Example 8** *The constraints given in example 2 are in fact generalized atomic constraints on the set (recall that we write LAdj for the 2 elements set  $\{\text{LAdj}_{belle}, \text{LAdj}_{ferme}\}$ ):*

$$\mathcal{U} = \{\text{Det}, \text{LAdj}, \text{RAdj}, \text{CN}, \text{Clit}, \text{TrV}, \text{IntrV}\}.$$

*Then the constraints are expressed on  $|\mathcal{U}| = 7$  elements and not on  $|\mathcal{G}| = 13$ .*

A generalized atomic constraint on  $\mathcal{U}$  can, of course, be expressed as a set of atomic constraints on  $\mathcal{G}$ : let  $u \in \mathcal{U}$  and  $t \in \mathcal{G}$  such that  $\bar{t} = u$

$$(\mathcal{L}, \mathcal{R}) \in \mathcal{C}(u) \implies \left( \bigcup_{L \in \mathcal{L}} L, \bigcup_{R \in \mathcal{R}} R \right) \in \mathcal{C}(t)$$

## 5.3 Implementation of lexicalized grammars

In implementations of large coverage linguistic resources, it is very common to have, first, the description of the set of “different” structures needed to describe the modeled natural language and then an anchoring mechanism that explains how words of the lexicon are linked to these structures. We call **unanchored grammar** the set  $\mathcal{U}$  of different structures (not yet related to words) that are needed to describe the grammar. In this context, the lexicon is split in two parts:

- a function  $\bar{\ell}$  from  $\mathcal{V}$  to subsets of  $\mathcal{U}$ ,
- an anchoring function  $\alpha$  which builds the grammar elements from a word  $w \in \mathcal{V}$  and an unanchored structure  $u \in \bar{\ell}(w)$ ; we suppose that  $\alpha$  verifies that  $\text{anc}(\alpha(w, u)) = w$ .

In applications, we suppose that  $\mathcal{U}$ ,  $\bar{\ell}$  and  $\alpha$  are given. In this context, we define the grammar as the codomain of the anchoring function:

$$\mathcal{G} = \bigcup_{w \in \mathcal{V}, u \in \bar{\ell}(w)} \alpha(w, u)$$

Now, we can define generalized constraints on the unanchored grammar, which are independent of the lexicon and can be computed statically for a given unanchored grammar.

## 6 Application to Interaction Grammars

In this section, we apply the Companionship Principle to the Interaction Grammars formalism. We first give a short and simplified description of IG and an example to illustrate them at work; we refer the reader to (Guillaume and Perrier, 2008) for a complete and detailed presentation.

### 6.1 Interaction Grammars

We illustrate some of the important features on the French sentence (2). In this sentence, “*la*” is an object clitic pronoun which is placed before the verb whereas the canonical place for the (non-clitic) object is on the right of the verb.

(2) “*Jean la demande.*” [John asks for it]

The set  $\mathcal{F}$  of final structures, used as output of the parsing process, contains ordered trees called **parse trees** (PT). An example of a PT for the sentence (2) is given in Figure 2. A PT for a sentence contains the words of the sentence or the empty word  $\epsilon$  in its leaves (the left-right order of the tree leaves follows the left-right order of words in the input sentence). The internal nodes of a PT represent the constituents of the sentence. The morpho-syntactic properties of these constituents are described with feature structures (only the category is shown in the figure).

As IG use the Model-Theoretic Syntax (MTS) framework, a PT is defined as the model of a set of constraints. Constraints are defined at the word level: words are associated to a set of constraints formally described as a **polarized tree description** (PTD). A PTD is a set of nodes provided with relations between these nodes. The three PTDs used to build the model above are given in Figure 3. The relations used in the PTDs are: immediate dominance (lines) and immediate sisterhood (arrows). Nodes represent syntactic constituents

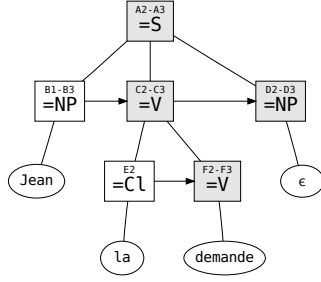


Figure 2: The PT of sentence (2)

and relations express structural dependencies between these constituents.

Moreover, nodes carry a polarity: the set of polarities is  $\{+, -, =, \sim\}$ . A  $+$  (resp.  $-$ ) polarity represents an available (resp. needed) resource, a  $\sim$  polarity describes a node which is unsaturated. Each  $+$  must be associated to exactly one  $-$  (and vice versa) and each  $\sim$  must be associated to at least another polarity.

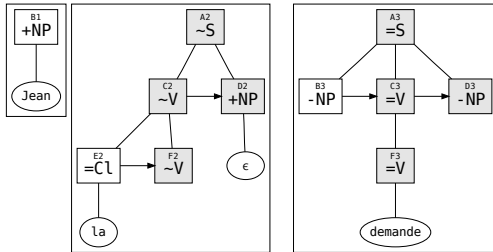


Figure 3: PTDs for the sentence (2)

Now, we define a PT to be a **model** of a set of PTDs if there is a surjective function  $\mathcal{J}$  from nodes of the PTDs to nodes of the PT such that:

- relations in the PTDs are realized in the PT: if  $M$  is a daughter (resp. immediate sister) of  $N$  in some PTD then  $\mathcal{J}(M)$  is a daughter (resp. immediate sister) of  $\mathcal{J}(N)$ ;
- each node  $N$  in the PT is saturated: the composition of the polarities of the nodes in  $\mathcal{J}^{-1}(N)$  with the associative and commutative rule given in Table 4 is  $=$ ;
- the feature structure of a node  $N$  in the PT is the unification of the feature structures of the nodes in  $\mathcal{J}^{-1}(N)$ .

One of the strong points of IG is the flexibility given by the MTS approach: PTDs can be partially superposed to produce the final tree (whereas superposition is limited in usual CG or in TAG for instance). In our example, the four grey nodes in the PTD which contains “*la*” are superposed to the four grey nodes in the PTD which contains “*demande*” to produce the four grey nodes in the model.

	$\sim$	$-$	$+$	$=$
$\sim$	$\sim$	$-$	$+$	$=$
$-$	$-$		$=$	
$+$	$+$	$=$		
$=$	$=$			

Figure 4: Polarity composition

In order to give an idea of the full IG system, we briefly give here the main differences between our presentation and the full system.

- Dominance relations can be underspecified: for instance a PTD can impose a node to be an ancestor of another one without constraining the length of the path in the model. This is mainly used to model unbounded extraction.
- Sisterhood relations can also be underspecified: when the order on subconstituents is not total, it can be modeled without using several PTDs.
- Polarities are attached to features rather than nodes: it sometimes gives more freedom to the grammar writer when the same constituent plays a role in different constructions.
- Feature values can be shared between several nodes: once again, this is a way to factorize the unanchored grammar.

The application of the Companionship Principle is described on the reduced IG but it can be straightforwardly extended to full IG with unessential technical details.

Following the notation given in 5.3, an IG is made of:

- A finite set  $\mathcal{V}$  of words;
- A finite set  $\mathcal{U}$  of unanchored PTDs (without any word attached to them);
- A lexicon function  $\bar{\ell}$  from  $\mathcal{V}$  to subsets of  $\mathcal{U}$ .



When  $t \in \bar{\ell}(w)$ , we can construct the anchored PTD  $\alpha(w, u)$ . Technically, in each unanchored PTD  $u$ , a place is marked to be the anchor, i.e. to be replaced by the word during the anchoring process. Moreover, the anchoring process can also be used to refine some features. The fact that the feature can be refined gives more flexibility and more compactness to the unanchored grammar construction. In the French IG grammar, the same unanchored PTD can be used for masculine or feminine common nouns and the gender is specified during the anchoring to produce distinct anchored PTDs for masculine and feminine nouns.  $\mathcal{G}$  is defined by:

$$\mathcal{G} = \bigcup_{w \in \mathcal{V}, u \in \bar{\ell}(w)} \alpha(w, u)$$

The parsing solutions of a lexical tagging is the set of PTs that are models of the list of PTDs described by the lexical tagging:

$$\mathbf{p}(L) = \{m \in \mathcal{F} \mid m \text{ is a model of } L\}$$

With the definitions of this section, an IG is a special case of AGF as defined in section 2.

## 6.2 Companionship Principle for IG

In order to apply the Companionship Principle, we have to explain how the generalized atomic constraints are built for a given grammar. One way is to look at dependency structures but in IG polarities are built in and then we can read the dependency information we need directly on polarities. A requirement to build a model is the saturation of all the polarities. This requirement can be expressed using atomic constraints. Each time a PTD contains an unsaturated polarity  $+$ ,  $-$  or  $\sim$ , we have to find some other compatible dual polarity somewhere else in the grammar to saturate it.

From the general MTS definition of IG above, we can define a step by step process to build models of a lexical tagging. The idea is to build incrementally the interpretation function  $\mathcal{J}$  with the atomic operation of **node merging**. In this atomic operation, we choose two nodes and make the hypothesis that they have the same image through  $\mathcal{J}$  and hence that they can be identified.

Now, suppose that the unanchored PTD  $u$  contains some unsaturated polarity  $p$ . We can use the atomic operation of node merging to test if the

unanchored PTD  $u'$  can be used to saturate the polarity  $p$ . Let  $\mathcal{L}$  (resp  $\mathcal{R}$ ) be the set of PTDs that can be used on the left (resp. on the right) of  $u$  to saturate  $p$ , then  $(\mathcal{L}, \mathcal{R})$  is a generalized atomic constraint in  $\mathcal{C}(u)$ .

## 7 Companionship Principle for other formalisms

As we said in the introduction, many current grammatical formalisms can more or less directly be used to generate dependency structures and hence are candidate for disambiguation with the Companionship Principle. With IG, we have seen that dependencies are strongly related to polarities: dependency constraints in IG are built with the polarity system.

We give below two short examples of polarity use to define atomic constraints on TAG and on CG. We use, as for IG, the polarity view of dependencies to describe how the constraints are built.

### 7.1 Tree Adjoining Grammars

Feature-based Tree Adjoining Grammars (hereafter FTAG) (Joshi, 1987) are a unification based version of Tree Adjoining Grammars. An FTAG consists of a set of elementary trees and of two tree composition operations: substitution and adjunction. There are two kinds of trees: auxiliary and initial. Substitution inserts a tree  $t$  with root  $r$  onto a leaf node  $l$  of another tree  $t'$  under the condition that  $l$  is marked as a place for substitution and  $l$  and  $r$  have compatible feature structures. Adjunction inserts an auxiliary tree  $t$  into a tree  $t'$  by splitting a node  $n$  of  $t'$  under the condition that the feature structures of the root and foot nodes of  $t$  are compatible with the *top* and *bottom* ones of  $n$ .

Getting the generalized atomic constraints and the model building procedure for lexical tagging is extremely similar to what was previously described for IG if we extend the polarization procedure which was described in (Gardent and Kow, 2005) to do polarity based filtering in FTAG. The idea is that for each initial tree  $t$ , its root of category  $C$  is marked as having the polarity  $+C$ , and its substitution nodes of category  $S$  are marked as having the polarity  $-S$ . A first constraint set contains trees  $t'$  whose root is polarized  $+S$  and such that feature structures are unifiable. A second constraint set contains trees  $t''$  which have a leaf that is polarized  $-C$ . We can extend this procedure to

auxiliary trees: each auxiliary tree  $t$  of category  $A$  needs to be inserted in a node of category  $A$  of another tree  $t'$ . This gives us a constraint in the spirit of the  $\sim$  polarity in IG:  $\mathcal{C}(t)$  contains all the trees  $t'$  in which  $t$  could be inserted<sup>2</sup>.

## 7.2 Categorical Grammars

In their type system, Categorical Grammars encode linearity constraints and dependencies between constituents. For example, a transitive verb is typed  $NP \setminus S / NP$ , meaning that it waits for a subject  $NP$  on its left and an object  $NP$  on its right. This type can be straightforwardly decomposed as two  $-NP$  and one  $+S$  polarities. Then again, getting the generalized atomic constraints is immediate and in the same spirit as what was described for IG.

## 7.3 Non-lexicalized formalisms

The lexicalization condition stated in section 2 excludes non-lexicalized formalisms like LFG or HPSG. Nothing actually prevents our method from being applied to these, but adding non-lexicalized combinators requires to complexify the formal account of the method. Adapting our method to HPSG would result in a generalization and unification of some of the techniques described in (Kiefer et al., 1999).

## 8 Experimental results

### 8.1 Setup

The experiments are performed using a French IG grammar on a set of 31 000 sentences taken from the newspaper *Le Monde*.

The French grammar we consider (Perrier, 2007) contains  $|\mathcal{U}| = 2\,088$  unanchored trees. It covers 88% of the grammatical sentences and rejects 85% of the ungrammatical ones on the TSNLP (Lehmann et al., 1996) corpus.

The constraints have been computed on the unanchored grammar as explained in section 5: each tree contains several polarities and therefore several atomic constraints. Overall, the grammar contains 20 627 atomic constraints. It takes 2 days to compute the set of constraints and the results can be stored in a constraints file of 10MB. Of course, an atomic constraint is more interesting when the sizes of  $\mathcal{L}$  and  $\mathcal{R}$  are small. In our grammar, 50% of the constraints set (either  $\mathcal{R}$  or  $\mathcal{L}$ )

<sup>2</sup>Note that in the adjunction case, the constraint is not oriented and then  $\mathcal{L} = \mathcal{R}$ .

contain at most 40 elements and 80% of these sets contain at most 200 elements over 2 088.

We give in figure 5 the number of sentences of each length in the corpus we consider.

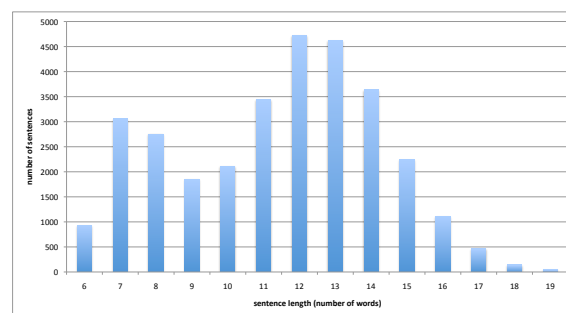


Figure 5: number of sentences of each length

## 8.2 Results

Two preliminary comments need to be made on the treatment of the results.

First, as we observed above, the number  $n$  of lexical taggings is a priori exponential in the length of the sentence. We thus consider its log. Moreover, because we use a raw corpus, some sentences are considered as ungrammatical by the grammar; in this case it may happen that the disambiguation method removes all taggings. In order to avoid undefined values when  $n = 0$ , we in fact consider  $\log_{10}(1 + n)$ .

Second, as expected, the ECP method is more time consuming and for some sentences the time and/or memory required is problematic. To be able to apply the ECP to a large number of sentences, we have used it after another filtering method based on polarities and described in (Bonfante et al., 2004).

Thus, for each sentence we have computed 3 different filters, each one finer than the previous:

- **QCP** the Quick Companionship Principle;
- **QCP+POL** QCP followed by a filtering technique based on polarity counting;
- **QCP+POL+ECP** the Exact Companionship Principle applied to the previous filter.

Figure 6 displays the mean computation time for each length: it confirms that the ECP is more time consuming and goes up to 5s for our long sentences.

Finally, we report the number of lexical taggings that each method returns. Figure 7 displays the mean value of  $\log_{10}(1 + n)$  where  $n$  is either the initial number of lexical taggings or the number of lexical taggings left by the filter.

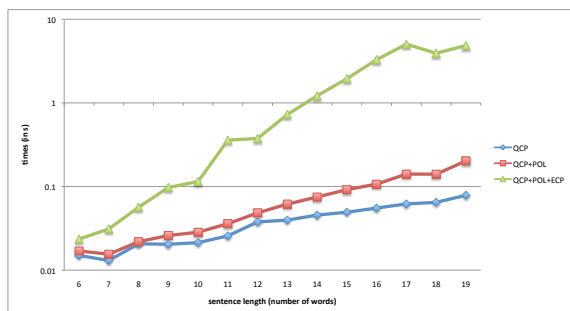


Figure 6: mean execution time (in s)

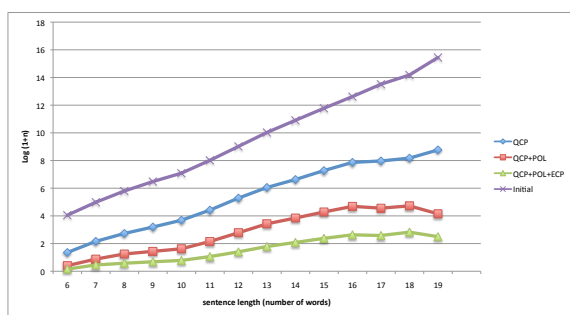


Figure 7: number of taggings (initial and after the 3 disambiguation methods)

We can observe that the slope of the lines corresponds to the mean word ambiguity: if the mean ambiguity is  $a$  then the number of taggings for a sentence of length  $n$  is about  $a^n$  and then  $\log(a^n) = n \cdot \log(a)$ . As a consequence, the mean ambiguity can be read as  $10^s$  where  $s$  is the slope in the last figure.

An illustration is given in figure 8 which exhibits the mean word ambiguity for sentences of length 16.

init	QCP	QCP+POL	QCP+POL+ECP
6.13	3.41	1.93	1.41

Figure 8: mean word ambiguity for sentences of length 16

## 9 Conclusion

We have presented a disambiguation method based on dependency constraints which allows to

filter out many wrong lexical taggings before entering the deep parsing. As this method relies on the computation of static constraints on the linguistic data and not on a statistical model, we can be sure that we will never remove any correct lexical tagging. Moreover, we manage to apply our methods to an interesting set of data and prove that it is efficient for a large coverage grammar and not only for a toy grammar.

These results are also an encouragement to develop further this kind of disambiguation methods. In the near future, we would like to explore some improvements.

First, we have seen that our principle cannot be computed on the whole grammar and that in its implementation we consider unanchored structures. We would like to explore the possibility of computing finer constraints (relative to the full grammar) on the fly for each sentence. We believe that this can eliminate some more taggings before entering the deep parsing.

Concerning the ECP, as we have seen, there is a kind of interplay between the efficiency of the filtering and the time of the computation. We would like to explore the possibility to define some intermediate way between QCP and ECP either by using approximate automata or using the ECP but only on a subset of elements where it is known to be efficient.

Another challenging method we would like to investigate is to use the Companionship Principle not only as a disambiguation method but as a guide for the deep parsing. Actually, we have observed for at least 20% of the words that dependencies are completely determined by the filtering methods. If deep parsing can be adapted to use this observation (this is the case for IG), this can be of great help.

Finally, we can improve the filtering using both worlds: the Companionship Principle and the polarity counting method. Two different constraints cannot be fulfilled by the same potential companion: this may allow to discover some more lexical taggings that can be safely removed.

## Acknowledgments

We would like to thank the anonymous reviewers for their helpful comments and suggestions.

## References

- Srinivas Bangalore and Aravind K. Joshi. 1999. Supertagging: an approach to almost parsing. *Comput. Linguist.*, 25(2):237–265.
- G. Bonfante, B. Guillaume, and G. Perrier. 2004. Polarization and abstraction of grammatical formalisms as methods for lexical disambiguation. In *CoLing 2004*, pages 303–309, Genève, Switzerland.
- P. Boullier. 2003. Supertagging : A non-statistical parsing-based approach. In *Proceedings of the 8th International Workshop on Parsing Technologies (IWPT 03)*, pages 55–65, Nancy, France.
- Stephen Clark and James R. Curran. 2004. The importance of supertagging for wide-coverage CCG parsing. In *COLING '04: Proceedings of the 20th international conference on Computational Linguistics*, page 282, Morristown, NJ, USA. Association for Computational Linguistics.
- S. Clark, J. Hockenmaier, and M. Steedman. 2002. Building Deep Dependency Structures with a Wide-Coverage CCG Parser. In *Proceedings of ACL'02*, pages 327–334, Philadelphia, PA.
- Ph. de Groote. 2001. Towards abstract categorial grammars. In *Association for Computational Linguistics, 39th Annual Meeting and 10th Conference of the European Chapter, Proceedings of the Conference*, pages 148–155.
- C. Gardent and E. Kow. 2005. Generating and selecting grammatical paraphrases. *Proceedings of the ENLG*, Aug.
- B. Guillaume and G. Perrier. 2008. Interaction Grammars. Research Report RR-6621, INRIA.
- A. Joshi and O. Rambow. 2003. A Formalism for Dependency Grammar Based on Tree Adjoining Grammar. In *Proceedings of the Conference on Meaning-Text Theory*.
- A. Joshi. 1987. An Introduction to Tree Adjoining Grammars. *Mathematics of Language*.
- S. Kahane. 2006. Polarized unification grammar. In *Proceedings of Coling-ACL'02*, Sydney.
- Bernd Kiefer, Hans-Ulrich Krieger, John Carroll, and Rob Malouf. 1999. A bag of useful techniques for efficient and robust parsing. In *Proceedings of the 37th annual meeting of the Association for Computational Linguistics on Computational Linguistics*, pages 473–480, Morristown, NJ, USA. Association for Computational Linguistics.
- A. Koller and M. Kuhlmann. 2009. Dependency trees and the strong generative capacity of CCG. In *EACL' 2009*, Athens, Greece.
- J. Kupiec. 1992. Robust Part-of-Speech Tagging Using a Hidden Markov Model. *Computer Speech and Language*, 6(3):225–242.
- F. Lamarche. 2008. Proof Nets for Intuitionistic Linear Logic: Essential Nets. Technical report, INRIA.
- J. Lambek. 1958. The mathematics of sentence structure. *American mathematical monthly*, pages 154–170.
- S. Lehmann, S. Oepen, S. Regnier-Prost, K. Netter, V. Lux, J. Klein, K. Falkedal, F. Fouvry, D. Estival, E. Dauphin, H. Compagnion, J. Baur, L. Balkan, and D. Arnold. 1996. TSNLP: Test Suites for Natural Language Processing. In *Proceedings of the 16th conference on Computational linguistics*, pages 711–716.
- J. Marchand, B. Guillaume, and G. Perrier. 2009. Analyse en dépendances à l'aide des grammaires d'interaction. In *Actes de TALN 09*, Senlis, France.
- B. Merialdo. 1994. Tagging English Text with a Probabilistic Model. *Computational linguistics*, 20:155–157.
- M. Moortgat and G. Morrill. 1991. Heads and phrases. Type calculus for dependency and constituent structure. In *Journal of Language, Logic and Information*.
- Takashi Ninomiya, Takuya Matsuzaki, Yoshimasa Tsuruoka, Yusuke Miyao, and Jun'ichi Tsujii. 2006. Extremely lexicalized models for accurate and fast HPSG parsing. In *Proceedings of the 2006 Conference on Empirical Methods in Natural Language Processing*, pages 155–163, Sydney, Australia, July. Association for Computational Linguistics.
- G. Perrier. 2007. A French Interaction Grammar. In *RANLP 2007*, pages 463–467, Borovets Bulgarie.
- M. Steedman. 2000. *The Syntactic Process*. MIT Press.
- L. Tesnière. 1959. *Éléments de syntaxe structurale*. Klinksiack.