

# Metagrammar Redux

Benoit Crabbé and Denys Duchier

LORIA, Nancy, France

**Abstract.** In this paper we introduce a general framework for describing the lexicon of a lexicalised grammar by means of elementary descriptive fragments. The system described hereafter consists of two main components: a control device aimed at controlling how fragments are to be combined together in order to describe meaningful lexical descriptions and a composition system aimed at resolving how elementary descriptions are to be combined.

## 1 Introduction

This paper is concerned with the design of large scaled grammars for natural language. It presents an alternative language of grammatical representation to the classical languages used for this purpose such as PATR II.

The need for a new language is motivated by the development of strongly lexicalised grammars based on tree structures rather than feature structures, and by the observation that, for tree based formalisms, lexical management with lexical rules raises non trivial practical issues [1].

In this paper we revisit a framework – the metagrammar – designed in particular for the lexical representation of tree based syntactic systems. It is articulated around two central ideas: (1) a *core* grammar is described by elementary tree fragments and (2) these fragments are combined by means of a control language to produce an *expanded* grammar. Throughout the paper, we illustrate the features of the framework using Tree Adjoining Grammar (TAG)[2] as a target formalism.

The paper is structured as follows. First (Section 2) we introduce the key ideas underlying grammatical representation taking PATR II as an illustration. We then provide the motivations underlying the design of our grammatical representation framework. The core metagrammatical intuition: lexical representation by manipulating fragments made of tree descriptions is provided in (Section 3). The motivations concerning the set up of an appropriate tree representation language are provided in Section 4. The fragment manipulation language is then developed in section 5. Section 6 introduces further questions concerning global conditions on model admissibility. And finally, the computational treatment of our description language is detailed in Section 7.

## 2 Lexical organisation

In this section we introduce the issue and the main ideas concerning lexical organisation of tree based syntactic systems. We begin by investigating the core

ideas developed in PATR II then we highlight inadequacies of PATR II for representing the lexicon of tree based syntactic systems such as Tree Adjoining Grammar.

*An historical overview: PATR II* Since the very first works in computational linguistics [3], lexical description roughly consists of specifying lexical entries together with a subcategorisation frame such as in PATR II:

```
love :
  <cat> = v
  <arg0 cat> = np
  <arg1 cat> = np
```

where we specify that the verb *love* takes two arguments: a *subject noun phrase* and an *object noun phrase*. This lexical entry, together with an appropriate grammar, is used to constrain the set of sentences in which *love* may be inserted. For instance this lexical entry is meant to express that *love* is used transitively as in *John loves mary* but not intransitively such as in *John loves* or *John loves to Mary*.

PATR II offers two devices to facilitate lexical description: templates and lexical rules. Templates are described by [3] as macros, and permit to easily state that that *love* and *write* are transitive verbs by writing:

```
love :
  transitiveVerb
write :
  transitiveVerb
transitiveVerb :
  <cat> = v
  <arg0 cat> = np
  <arg1 cat> = np
```

where `transitiveVerb` is a macro called in the descriptions of *love* and *write*. On the other hand, lexical rules are used to describe multiple variants of verbs. For instance, to express that a transitive verb such as *love* may be used in its active or passive variant we may add the following lexical rule to our lexicon:

```
passive :
  <out cat> = <in cat>
  <out arg1 cat> = <in arg0 cat>
  <out arg0 cat> = pp
```

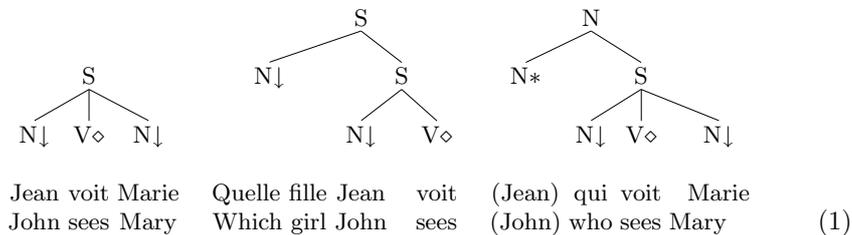
This rule says that a new lexical entry `out` is to be build from an initial lexical entry `in` where the category of `out` is identical to the category of `in`, the category of the object becomes the category of the subject and that the subject category now becomes prepositional phrase.

Lexical rules are meant to allow a dynamic expansion of related lexical variants. So for the verb *love* the application of the `passive` lexical rule to its base entry generates a new, derived, passive lexical entry meaning that both active and passive variants are licensed by the lexical entries.

Variants and improvements of this classical system have been (and are still) used for describing the lexicon in various syntactic frameworks such as LFG[4] or HPSG [5]. Whatever the differences, two leading ideas remain nowadays: lexical description aims both at factorising information (templates) and at expressing relationships between variants of a same lexical unit (lexical rules).

*Tree Adjoining Grammar: a case study* Tree adjoining grammar (TAG)<sup>1</sup> is a tree composition system aimed at describing natural language syntax [2] which strongly lexicalised. In other words, a tree adjoining grammar consists of a lexicon, the elementary trees, each of them being associated to a lexical unit, and two operations used for combining the lexical units: adjunction and substitution.

Following the conventions used in TAG implementations such as XTAG [6], we work with tree schematas (or templates) such as these<sup>2</sup>:



where the appropriate lexical word is inserted dynamically by the parser as a child of the anchor (marked  $\diamond$ ). The nodes depicted with an arrow ( $\downarrow$ ) are the substitution nodes and those depicted with a star ( $\star$ ) are the foot nodes.

Strikingly, works concerning lexical organisation of strongly lexicalised syntactic systems often try to provide alternative solutions to that of Shieber. The main reason is that the amount and the variety of lexical units is much more important, therefore the number of templates and lexical rules to be used is strongly increased. In the context of the development of large sized grammars, this situation requires the grammar writer to design complicated ordering schemes as it is illustrated by [1].

To overcome this, we take up an idea first introduced in [8] for Construction Grammar. Roughly speaking they describe the lexicon using a dynamic process: given a *core* lexicon manually described they build up an *expanded* lexicon by combining elementary *fragments* of information.

Besides strong lexicalisation, setting up a system representing a TAG lexicon raises another problem, that of the structures used. In Construction Grammar, [8] combine elementary fragments of information via feature structure unification. When working with TAG, however, one works with trees

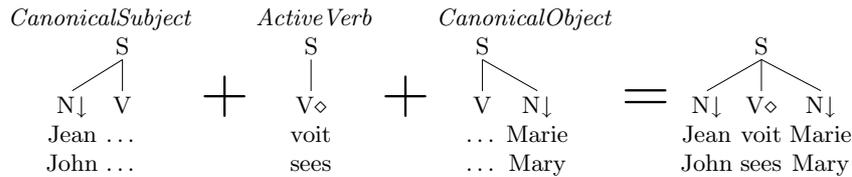
<sup>1</sup> Strictly speaking, we mean here Lexicalised Tree Adjoining Grammar (LTAG). Indeed, the system is usually used in its lexicalised version[6].

<sup>2</sup> The trees depicted in this paper are motivated by the French grammar of [7] who provides linguistic justifications in particular for the non utilisation of the VP category and the use of N at multiple bar levels instead of introducing the category NP in French.

### 3 Introduction to the framework

In this section we sketch the idea of describing the lexicon by controlling combinations of elementary fragment descriptions.

This idea stems from the following observation: the design of a TAG grammar consists of describing trees made of elementary pieces of information (hereafter: fragments). For instance the following tree is defined by combining a subtree representing a subject another subtree representing an object and finally a subtree representing the spine of the verbal tree:



Of course, we will also want convenient means of expressing variants of the above tree; for example, where, the subject instead of being realized in canonical position is realized as a questioned subject (*wh*) or a relative subject.

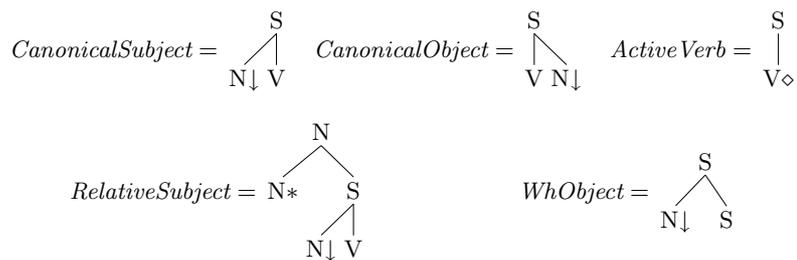
More generally while designing a grammar one wants to express general statements for describing sets of trees: for instance, a transitive verb is made of a subject, an object and a verbal active spine. In short we would like to write something like:

$$\text{TransitiveVerb} = \text{Subject} \wedge \text{ActiveVerb} \wedge \text{Object}$$

where *Subject* and *Object* are shortcuts for describing sets of variants:

$$\begin{aligned}
 \text{Subject} &= \text{CanonicalSubject} \vee \text{RelativeSubject} \\
 \text{Object} &= \text{CanonicalObject} \vee \text{WhObject}
 \end{aligned}$$

and where *CanonicalSubject*, *WhSubject*... are defined as the actual fragments of the grammar:



Given the above definitions, a description such as *TransitiveVerb* is intended to describe the following tree schematas depicted in (1)<sup>3</sup>. That is each variant

<sup>3</sup> The combination of relative subject and a questioned object is rejected by the principle of extraction unicity (See section 6).

description of the subject embedded in the *Subject* clause is combined with each variant description of the object in the *Object* clause and the description in the *Active Verb* clause.

As it stands, the representation system we have introduced so far requires to set up two components: first we investigate which language to use for describing *tree fragments* and combining them (Section 4). Second we detail the language which controls how fragments are to be combined (Section 5).

## 4 A language for describing tree fragments

In this section, we consider two questions: (1) how to conveniently describe tree fragments, (2) how to flexibly constrain how such tree fragments maybe combined to form larger syntactic units. We first introduce a language of tree descriptions, and then show how it can be generalized to a family of formal languages parametrized by an arbitrary constraining decoration system that further limits how elements can be combined.

*The base language L.* Let  $x, y, z \dots$  be nodes variables. We write  $\triangleleft$  for immediate dominance,  $\triangleleft^*$  for its reflexive transitive closure (dominance),  $\prec$  for immediate precedence (or adjacency) and  $\prec^+$  for its transitive closure (strict precedence). We let  $\ell$  range over a set of node labels generally intended to capture the notion of categories. A tree description  $\phi$  has the following abstract syntax:

$$\phi ::= x \triangleleft y \mid x \triangleleft^* y \mid x \prec y \mid x \prec^+ y \mid x : \ell \mid \phi \wedge \phi \quad (2)$$

$L$ -descriptions are, as expected, interpreted over first-order structures of finite, ordered, constructor trees. As usual, we limit our attention to minimal models.

Throughout the paper we use an intuitive graphical notation for representing tree descriptions. Though this notation is not sufficient to represent every expression of the language, it nonetheless generally suffices for the kind of trees typically used in natural language syntax. Thus, the description  $(D_0)$  on the left is graphically represented by the tree notation on the right:

$$\begin{array}{l} x \triangleleft^* w \wedge x \triangleleft y \wedge x \triangleleft z \\ D_0 = \wedge y \prec^+ z \wedge z \prec w \\ \wedge x : X \wedge y : Y \wedge z : Z \wedge w : W \end{array} \quad (D_0) \quad \begin{array}{c} X \\ \swarrow \quad \downarrow \quad \searrow \\ Y \prec^+ Z \quad W \end{array} \quad (3)$$

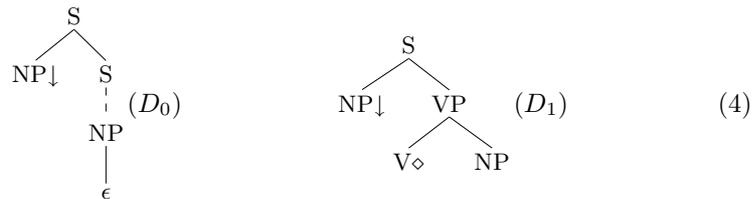
where immediate dominance is represented by a solid line, dominance by a dashed line, precedence by the symbol  $\prec^+$  and adjacency is left unmarked.

*A parametric family of languages.* It is possible to more flexibly control how tree fragments maybe combined by adding annotations to nodes together with stipulations for how these annotations restrict admissible models and interpretations. In this manner, we arrive at the idea of a family of languages  $L(C)$  parametrized by a *combination schema C*.

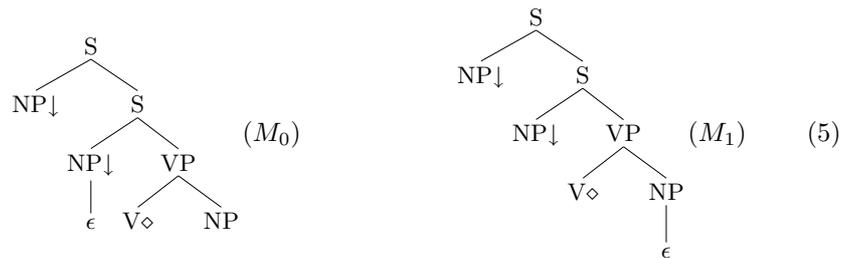
In the remainder of this section we discuss three instantiations of  $L(C)$  that have been used for describing the lexicon of Tree Adjoining Grammars. The first

one,  $L(\emptyset)$  is used by Xia [9], the second one  $L(names)$  is used by Candito [10]. We show that neither  $L(\emptyset)$  nor  $L(names)$  are appropriate for describing the lexicon of a French TAG Grammar. We then introduce  $L(colors)$  which we have used successfully for that purpose.

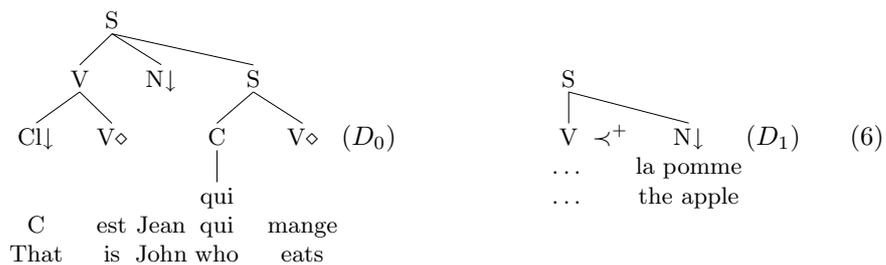
*Language  $L(\emptyset)$ .* This first instantiation of  $L(C)$  is used by F. Xia[9]. This language does not use any combination constraint. The combination schema  $C$  is thus empty. Equipped with such a language we can independently describe fragments such as these<sup>4</sup>:



where  $D_0$  describes a relative NP and  $D_1$  a transitive construction. Their combination leads to the following two models:

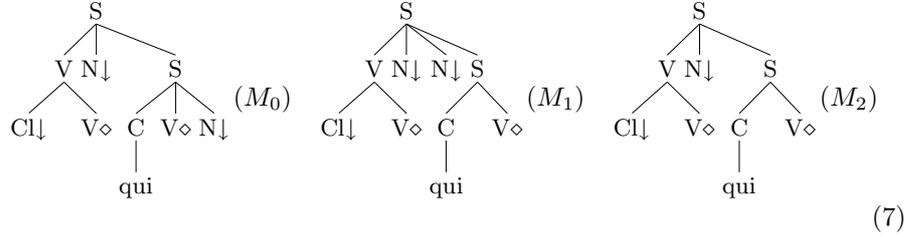


However this language faces an expressivity limit since, for the purpose of lexical organisation, linguists want to constrain combinations more precisely. For instance, in the French Grammar the following fragment composition is badly handled since:



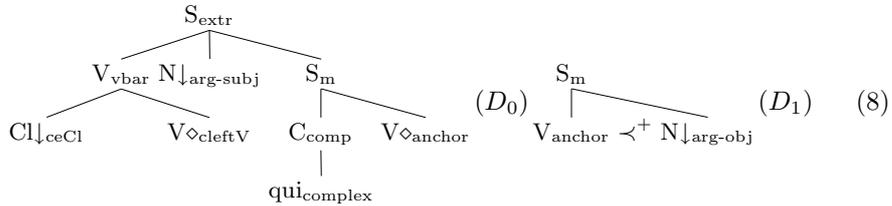
<sup>4</sup> These fragments and the related models are those used by F. Xia in the context of the XTAG English Grammar.

yields, among others, the following results:



where only  $M_0$  is normally deemed linguistically valid.

*Language  $L(\text{names})$ .* In her thesis, M.-H. Candito [10] introduces an instance of  $L(C)$  that constrains combinations to avoid cases such as the one outlined above. The combination schema  $C$  is as follows: (1) a finite set of names where each node of a tree description is associated to such a name and (2) Two nodes sharing the same name are to be interpreted as denoting the same entity, hence when merging descriptions, only the nodes with the same names are merged. In other words, a model is valid if (1) every node has exactly one name and (2) there is at most one node with a given name<sup>5</sup>.



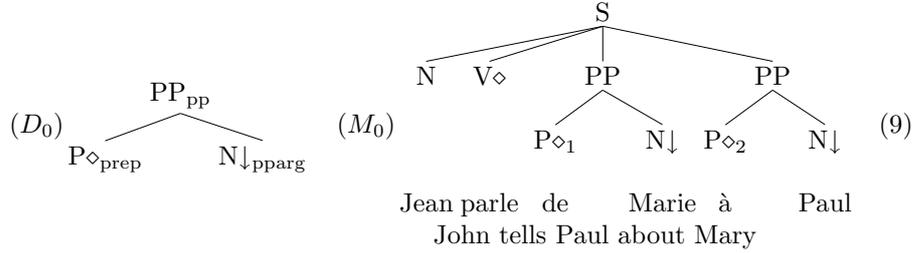
The model resulting from merging  $D_0$  with  $D_1$  is only  $M_0$  depicted in (7). In such a case,  $L(\text{names})$  corrects the shortcomings of  $L(\emptyset)$ . However, during the development of a non trivial grammar using this language, it turned out that  $L(\text{names})$  was eventually unsatisfactory for two main reasons:

The first is practical and rather obvious: the grammar writer has to manage naming by hand, and must handle the issues arising from name collisions.

The second is more tricky: the grammar writer may need to use the same tree fragment more than once in the same description. For example, such an

<sup>5</sup> To be complete, M.-H. Candito uses additional operations to map multiple names on a single node. However this does not change the content of our actual discussion.

occasion arises in the case of a double PP complementation:



where one cannot use the fragment  $(D_0)$  more than once to yield  $M_0$  since identical names must denote identically the same nodes.

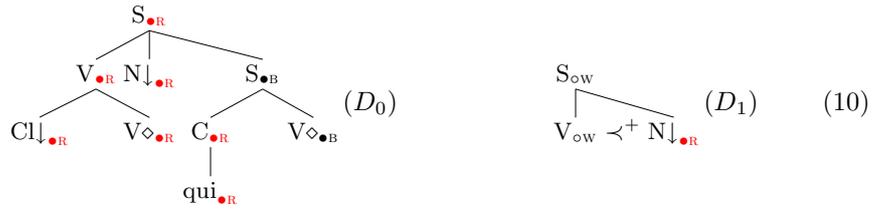
*A Language with colored nodes  $L(colors)$ .* We used this language in the development of a large scale French TAG patterned after the analysis of [7].

$L(colors)$  was designed to overcome the shortcomings of languages  $L(\emptyset)$  and  $L(names)$ . We want (1) to be able to constrain more precisely the way fragments combine together than with language  $L(\emptyset)$  (2) we want to eschew the explicit naming management of language  $L(names)$ .

To do this, the combination schema  $C$  used in  $L(colors)$  decorates all nodes with colors: black ( $\bullet_B$ ), white ( $\circ_W$ ), red ( $\bullet_R$ ) or failure ( $\perp$ ). The additional condition on model admissibility is that each node must be either red or black.

When combining tree descriptions, nodes are merged and their colors combined. The table to the right specifies the result of combining two colors. For instance, combining a white node with a black node yields a black node; combining a white node with a red node is illegal and produces a failure. As a matter of illustration, the following color enriched descriptions yields only the desired model  $(M_0)$  for example number (7)<sup>6</sup>

	$\bullet_B$	$\bullet_R$	$\circ_W$	$\perp$
$\bullet_B$	$\perp$	$\perp$	$\bullet_B$	$\perp$
$\bullet_R$	$\perp$	$\perp$	$\perp$	$\perp$
$\circ_W$	$\bullet_B$	$\perp$	$\circ_W$	$\perp$
$\perp$	$\perp$	$\perp$	$\perp$	$\perp$



Intuitively the colors have a semantic similar to that of resources and requirements systems such as Interaction Grammars [11]. A tree is well formed if it is saturated. The colors representing saturation are red or black the color representing non saturation is white and we have a color representing failure.

<sup>6</sup> We let the reader figure out how to express double PP complementation (9). It requires to use a description similar to  $(D_1)$  depicted here, patterned for describing a prepositional phrase though.

Though  $L(\text{colors})$  turned out to be satisfactory for designing a large scale French TAG, it might not be similarly adequate for other frameworks or languages.<sup>7</sup> However, alternative instances of  $L(C)$  might be suitable. For example a combination schema based on polarities seems a very reasonable foundation for interaction grammars [11] and even for polarity based unification grammars [12].

## 5 Controlling fragment combinations

In Section 3 we identified a number of desirable requirements for a metagrammar language: (1) it should support disjunctions to make it easy to express diathesis (such as active, passive), (2) it should support conjunction so that complex descriptions can be assembled by combining several simpler ones, (3) it should support abstraction so that expressions can be named to facilitate reuse and avoid redundancy.

In this section, we introduce the language  $\mathcal{L}_C$  to control how fragments can be combined in our proposed lexical representation framework, and show how  $\mathcal{L}_C$  satisfies all the above requirements.

$$\text{Clause} ::= \text{Name} \rightarrow \text{Goal} \quad (11)$$

$$\text{Goal} ::= \text{Goal} \wedge \text{Goal} \quad | \quad \text{Goal} \vee \text{Goal} \quad | \quad \phi \quad | \quad \text{Name} \quad (12)$$

This language allows to manipulate fragment descriptions ( $\phi$ ), to express the composition of statements ( $\text{Goal} \wedge \text{Goal}$ ), to express nondeterministic choices ( $\text{Goal} \vee \text{Goal}$ ), and finally to name complex statements for reuse ( $\text{Name} \rightarrow \text{Goal}$ ).

The main motivation for the control language is to support the combination and reuse of tree fragments. Instead of manipulating directly tree descriptions, the language allows to define abstractions over (possibly complex) statements. Thus, the clause:

$$\text{CanonicalSubject} \rightarrow \begin{array}{c} \text{S} \\ \swarrow \downarrow \\ \text{N} \downarrow \text{V} \end{array} \quad (13)$$

defines the abstraction *CanonicalSubject* to stand for a tree description which can be subsequently reused via this new name, while the clause:

$$\text{TransitiveVerbActive} \rightarrow \text{Subject} \wedge \text{ActiveVerb} \wedge \text{Object} \quad (14)$$

states that a lexical tree for a transitive verb is formed from the composition of the descriptions of a subject, of an object and of an active verb.

Disjunction is interpreted as an nondeterministic choice: each of the alternatives describes one of the ways in which the abstraction can be realized. As

---

<sup>7</sup> The current framework is not restricted to the specific case of Tree Adjoining Grammars. It should be straightforward to adapt it to other cases of tree based syntactic systems such as Interaction Grammars.

illustrated by lexical rules as used e.g. in PATR II [3], a system of lexical representation needs to be equipped with a way to express relationships between lexical items such as does a passive lexical rule relating an active and a passive lexical entry. In our approach, similar relations are expressed with disjunctions. Thus the following statement expresses the fact that various realisation of the subject are equivalent:

$$\begin{aligned}
 \textit{Subject} &\rightarrow \textit{CanonicalSubject} & (15) \\
 &\vee \textit{WhSubject} \\
 &\vee \textit{RelativeSubject} \\
 &\vee \textit{CliticSubject}
 \end{aligned}$$

As surely has become evident, the language presented in this section has very much the flavor of a logic programming language. More precisely, it can be understood as an instance of the *Definite Clause Grammar* (DCG) paradigm. DCGs were originally conceived to express the production rules of context free grammars: they characterized the sentences of a language, i.e. all the possible ways words could be combined into grammatical sequences by concatenation. Here, instead of words, we have tree fragments, and instead of concatenation we have a composition operation. In other words,  $\mathcal{L}_C$  allows us to write the grammar of a tree grammar, which surely justifies the name *metagrammar*.

## 6 Principles of well-formedness

So far, the current system assumes that one can describe grammatical information by combining fragments of local information. There are however cases where the local fragments interact when realised together. To handle these interactions in an elegant way, the system allows to formulate additional global constraints on tree admissibility, called the principles.

Let us express in the control language the fact that a transitive verb is made of a subject, an object and a verb in the active form:

$$\textit{TransitiveVerb} \rightarrow \textit{Subject} \wedge \textit{ActiveVerb} \wedge \textit{Object} \quad (16)$$

$$\textit{Subject} \rightarrow \textit{CanonicalSubject} \vee \textit{CliticSubject} \quad (17)$$

$$\textit{Object} \rightarrow \textit{CanonicalObject} \vee \textit{CliticObject} \quad (18)$$

*Clitic ordering* According to the subject and object clauses, it is the case that among others, a description of a transitive verb is made of the composition of a clitic subject and a clitic object<sup>8</sup> whose definitions are as follows:

$$\begin{array}{ccc}
 \textit{CliticSubject} \rightarrow & \begin{array}{c} \text{V} \\ \swarrow \quad \searrow \\ \text{Cl}\downarrow[\text{case} = \text{nom}] \quad \leftarrow^+ \text{V} \end{array} & \textit{CliticObject} \rightarrow \begin{array}{c} \text{V} \\ \swarrow \quad \searrow \\ \text{Cl}\downarrow[\text{case} = \text{acc}] \quad \leftarrow^+ \text{V} \end{array} & (19)
 \end{array}$$

<sup>8</sup> In French, clitics are small non tonic pronominal particles realized in front of the verb which are ordered according to a fixed order. The problem of clitic ordering is a well known case of such an interaction. It was already described as problematic in the Generative literature in the early 70's [13].

When realized together, none of the clitic descriptions say how these clitics are ordered relative to each other so that a merge of these two descriptions yields the following two models:



where  $M_1$  is an undesirable solution in French.

French clitic ordering is thus rendered by a principle of tree well formedness: *sibling nodes of category Clitic have to be ordered according to the respective order of their ranking property*. So, if we take the case feature of descriptions (19) to be the ranking property, and that the order defined over the property constrains (inter alia) nominative to precede accusative then in every tree where both a nominative and an accusative clitic are realised, the principle ensures that only  $M_0$  is a valid model.

*Extraction unicity* Another principle presented hereafter (Section 7) is that of extraction unicity. We assume that, in French, only one argument of a given predicate may be extracted<sup>9</sup>. Following this, the extraction principle is responsible for ruling out trees models where more than a node would be associated to the property of extraction.

Two other principles have actually been used in the implementation of the French Grammar: a principle for ensuring clitic unicity and a principle for expressing islands constraints<sup>10</sup>. The expression of an additional principle of functional unicity is currently under investigation.

## 7 A constraint satisfaction approach

As mentioned earlier, the control language  $\mathcal{L}_C$  of Section 5 can be regarded as an instance of the *Definite Clause Grammar* (DCG) paradigm. While DCGs are most often used to describe sentences, i.e. sequences of words, here, we apply them to the description of formulae in language  $L(\text{colors})$ , i.e. conjunctions of colored tree fragments.

A consequence of regarding a metagrammar, i.e. a program expressed in language  $\mathcal{L}_C$ , as a DCG is that it can be reduced to a logic program and executed as such using well-known techniques. What remains to be explained is how, from a conjunction of colored tree fragments, we derive all complete trees that can be formed by combining these fragments together.

For this task, we propose a constraint-based approach that builds upon and extends the treatment of dominance constraints of Duchier and Niehren [15]. We begin by generalizing slightly the language introduced in Section 4 to make

<sup>9</sup> Actually, cases of double extraction have been discovered in French, they are so rare and so unnatural that they are generally ruled out of the grammatical implementations.

<sup>10</sup> This principle is related to the way one formalises islands constraints in TAG [14].

it more directly amenable to the treatment described in [15], then we show how we can enumerate the minimal models of a description in that language by translating this description into a system of constraints involving set variables, and solving that instead.

*Tree description language.* In order to account for the idea that each node of a description is colored either red, black or white, we let  $x, y, z$  range over 3 disjoint sets of node variables:  $V_R, V_B, V_W$ . We write  $\triangleleft$  for immediate dominance,  $\triangleleft^+$  for its transitive closure, i.e. strict dominance,  $\prec$  for immediate precedence, and  $\prec^+$  for its transitive closure, i.e. strict precedence. We let  $\ell$  range over a set of node labels. A description  $\phi$  has the following abstract syntax:

$$\phi ::= x R y \mid x \triangleleft y \mid x \prec y \mid x : \ell \mid \phi \wedge \phi \quad (20)$$

where  $R \subseteq \{=, \triangleleft^+, \triangleright^+, \prec^+, \succ^+\}$  is a set of relation symbols whose intended interpretation is disjunctive; thus  $x \{=, \triangleleft^+\} y$  is more conventionally written  $x \triangleleft^* y$ .

In [15], the abstract syntax permitted a literal of the form  $x : \ell(x_1, \dots, x_n)$  that combined (1) an assignment of the label  $\ell$  to  $x$ , (2) immediate dominance literals  $x \triangleleft x_i$ , (3) immediate precedence literals  $x_i \prec x_{i+1}$ , (4) an arity constraint stipulating that  $x$  has exactly  $n$  children. Here we prefer a finer granularity and admit literals for immediate dominance and immediate precedence. For simplicity of presentation we omit an arity constraint literal.

*Enumerating minimal models.* We now describe how to convert a description into a constraint system that uses set constraints and such that the solutions of the latter are in bijection with the minimal models of the former. Such a constraint system can be realized and solved efficiently using the constraint programming support of Mozart/Oz. Our conversion follows very closely the presentation of [15].

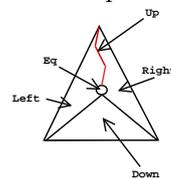
The general intuition is that a literal  $x R y$  should be represented by a membership constraint  $y \in R(x)$  where  $R(x)$  is a set variable denoting all the nodes that stand in  $R$  relationship with  $x$ . We write  $V^\phi$  for the set of variables occurring in  $\phi$ . Our encoding consists of 3 parts:

$$\llbracket \phi \rrbracket = \bigwedge_{x \in V^\phi} A_1(x) \wedge \bigwedge_{x, y \in V^\phi} A_2(x, y) \wedge B[\phi] \quad (21)$$

$A_1(\cdot)$  introduces a node representation per variable,  $A_2(\cdot, \cdot)$  axiomatizes the tree-ness of the relations between these nodes, and  $B(\cdot)$  encodes the problem-specific restrictions imposed by  $\phi$ .

## 7.1 Representation

When observed from a specific node  $x$ , the nodes of a solution tree (a model), and hence the variables which they interpret, are partitioned into 5 regions: the node denoted by  $x$  itself, all nodes below, all



nodes above, all nodes to the left, and all nodes to the right. The main idea is to introduce corresponding set variables  $Eq_x, Up_x, Down_x, Left_x, Right_x$  to encode the sets of variables that are interpreted by nodes in the model which are respectively equal, above, below, left, and right of the node interpreting  $x$ . First, we state that  $x$  is one of the variables interpreted by the corresponding node in the model:

$$x \in Eq_x \quad (22)$$

Then, as explained above, we have the following fundamental partition equation:

$$V^\phi = Eq_x \uplus Up_x \uplus Down_x \uplus Left_x \uplus Right_x \quad (23)$$

where  $\uplus$  denotes *disjoint union*. We can (and in fact *must*, as proven in [15]) improve propagation by introducing shared intermediate results  $Side, Eqdown, Equip, Eqdownleft, Eqdownright$ .

$$Side_x = Left_x \uplus Right_x \quad Eqdownleft_x = Eqdown_x \uplus Left_x \quad (24)$$

$$Eqdown_x = Eq_x \uplus Down_x \quad Eqdownright_x = Eqdown_x \uplus Right_x \quad (25)$$

$$Equip_x = Eq_x \uplus Up_x \quad (26)$$

which must all be related to  $V^\phi$ :

$$V^\phi = Eqdown_x \uplus Up_x \uplus Side_x \quad V^\phi = Eqdownleft_x \uplus Up_x \uplus Right_x \quad (27)$$

$$V^\phi = Equip_x \uplus Down_x \uplus Side_x \quad V^\phi = Eqdownright_x \uplus Down_x \uplus Left_x \quad (28)$$

We define  $A_1(x)$  as the conjunction of the constraints introduced above.

## 7.2 Wellformedness

Posing  $\mathbf{Rel} = \{=, \triangleleft^+, \triangleright^+, \triangleleft^-, \triangleright^-\}$ , in a tree, the relationship that obtains between the nodes denoted by  $x$  and  $y$  must be one in  $\mathbf{Rel}$ : the options are mutually exclusive. We introduce a variable  $C_{xy}$ , called a *choice variable*, to explicitly represent it and contribute a well-formedness clause  $A_3[x r y]$  for each  $r \in \mathbf{Rel}$ .

$$A_2(x, y) = C_{xy} \in \mathbf{Rel} \wedge \wedge \{A_3[x r y] \mid r \in \mathbf{Rel}\} \quad (29)$$

$$A_3[x r y] \equiv D[x r y] \wedge C_{xy} = r \vee C_{xy} \neq r \wedge D[x \neg r y] \quad (30)$$

For each  $r \in \mathbf{Rel}$ , it remains to define  $D[x r y]$  and  $D[x \neg r y]$  encoding respectively the relationships  $x r y$  and  $x \neg r y$  by set constraints on the representations of  $x$  and  $y$ .

$$D[x = y] = Eq_x = Eq_y \wedge Up_x = Up_y \wedge \dots \quad (31)$$

$$D[x \neg = y] = Eq_x \parallel Eq_y \quad (32)$$

$$D[x \triangleleft^+ y] = Eqdown_y \subseteq Down_x \wedge Equip_x \subseteq Up_y \wedge Left_x \subseteq Left_y \wedge Right_x \subseteq Right_y \quad (33)$$

$$D[x \neg \triangleleft^+ y] = Eq_x \parallel Up_y \wedge Down_x \parallel Eq_y \quad (34)$$

$$D[x \triangleleft^- y] = Eqdownleft_x \subseteq Left_y \wedge Eqdownright_y \subseteq Right_x \quad (35)$$

$$D[x \neg \triangleleft^- y] = Eq_x \parallel Left_y \wedge Right_x \parallel Eq_y \quad (36)$$

where  $\parallel$  represents disjointness.

### 7.3 Problem-specific constraints

The third part  $\mathbb{B}[\phi]$  of the translation forms the problem-specific constraints that further restrict the admissibility of well-formed solutions and only accepts those which are models of  $\phi$ . The translation is given by case analysis following the abstract syntax of  $\phi$ :

$$\mathbb{B}[\phi \wedge \phi'] = \mathbb{B}[\phi] \wedge \mathbb{B}[\phi'] \quad (37)$$

A rather nice consequence of introducing choice variables  $C_{xy}$  is that any dominance constraint  $x R y$  can be translated as a restriction on the possible values of  $C_{xy}$ . For example  $x \triangleleft^* y$  can be encoded as  $C_{xy} \in \{=, \triangleleft^+\}$ . More generally:

$$\mathbb{B}[x R y] = C_{xy} \in R \quad (38)$$

A labeling literal  $x : \ell$  simply restricts the label associated with variable  $x$ :

$$\mathbb{B}[x : \ell] = \text{Label}_x = \ell \quad (39)$$

An immediate dominance literal  $x \triangleleft y$  not only states that  $x \triangleleft^+ y$  but also that there are no intervening nodes on the spine that connects the two nodes:

$$\mathbb{B}[x \triangleleft y] = C_{xy} = \triangleleft^+ \wedge \text{Up}_y = \text{Equip}_x \quad (40)$$

An immediate precedence literal  $x \prec y$  not only states that  $x \prec^+ y$  but also that there are no intervening nodes horizontally between them:

$$\mathbb{B}[x \prec y] = C_{xy} = \prec^+ \wedge \text{Eqdownleft}_x = \text{Left}_y \wedge \text{Right}_x = \text{Eqdownright}_y \quad (41)$$

### 7.4 Well-coloring

While distinguishing left and right is a small incremental improvement over [15], the treatment of colors is a rather more interesting extension. The main question is: which nodes can or must be identified with which other nodes? Red nodes cannot be identified with any other nodes. Black nodes may be identified with white nodes. Each white node must be identified with a black node. As a consequence, for every node, there is a unique red or black node with which it is identified. We introduce the (integer) variable  $RB_x$  to denote the red or black node with which  $x$  is identified.

For a red node,  $x$  is identified only with itself:

$$x \in V_R \quad \Rightarrow \quad RB_x = x \wedge \text{Eq}_x = \{x\} \quad (42)$$

For a black node, the constraint is a little relaxed (it may also be identified with white nodes):

$$x \in V_B \quad \Rightarrow \quad RB_x = x \quad (43)$$

Posing  $V_B^\phi = V^\phi \cap V_B$ , each white node must be identified with a black node:

$$x \in V_W \quad \Rightarrow \quad RB_x \in V_B^\phi \quad (44)$$

Additionally, it is necessary to ensure that  $RB_x = RB_y$  iff  $x$  and  $y$  have been identified. We can achieve this simply by modifying the definition (32) of  $D[x \neg= y]$  as follows:

$$D[x \neg= y] \quad = \quad Eq_x \parallel Eq_y \wedge RB_x \neq RB_y \quad (45)$$

## 7.5 Extraction Principle

As an illustration of how the framework presented so far can be extended with linguistically motivated principles to further constrain the admissible models, we describe now what we have dubbed the *extraction principle*.

The description language is (somehow) extended to make it possible to mark certain nodes of a description as representing an *extraction*. The extraction principle then makes the additional stipulation that, to be admissible, a model must contain at most one node marked as extracted.

Let  $V_{\text{XTR}}^\phi$  be the subset of  $V^\phi$  of those node variables marked as extracted. We introduce the new boolean variable  $Extracted_x$  to indicate whether the node denoted by  $x$  is extracted:

$$Extracted_x \quad = \quad Eq_x \cap V_{\text{XTR}}^\phi \neq \emptyset \quad (46)$$

Posing  $V_{\text{RB}}^\phi = V^\phi \cap (V_R \cup V_B)$ , and freely identifying the boolean values false and true respectively with the integers 0 and 1, the extraction principle can be enforced with the following constraint:

$$\sum_{x \in V_{\text{RB}}^\phi} Extracted_x \quad < \quad 2 \quad (47)$$

## 8 Conclusion

This paper introduces a core abstract framework for representing grammatical information of tree based syntactic systems. Grammatical representation is organised around two central ideas: (1) the lexicon is described by means of elementary tree fragments that can be combined. (2) Fragment combinations are handled by a control language, which turns out to be an instance of a DCG.

The framework described here, generalises the TAG specific approaches of [9, 10]. We have provided a parametric family of languages for tree composition as well as constraints on tree well formedness.

Besides the non TAG specific tree composition language, it mostly differs from the TAG instantiations by (1) it introduces a control language allowing to express explicitly composition of fragments as well as variants of related lexical

entries. The two existing systems of [10] and [9] rely mostly on an algorithmic device for expressing variants, namely a crossing algorithm for [10], and an external module of lexical rules for [9].

The introduction of the control language (1) avoids to work with different modules and (2) introduces more flexibility in expressing variants that avoids to deal with "shadow" classes as it turns out to be the case in [10].

The framework presented here has been extensively tested against the development of a large sized French TAG based on [7]. This grammar covers most of the phenomenons related to the syntax of French verbs.

## References

1. Prolo, C.: Generating the xtag english grammar using metarules. In: Proc. COLING 2002, Taiwan (2002)
2. Joshi, A.K., Schabès, Y.: Tree adjoining grammars. In Rozenberg, G., Salomaa, A., eds.: Handbook of Formal Languages. Springer Verlag, Berlin (1997)
3. Shieber, S.M.: The design of a computer language for linguistic information. In: Proceedings of the Tenth International Conference on Computational Linguistics, Stanford University, Stanford, California (1984) 362–366
4. Kaplan, R.M., Maxwell, J.T.: Lfg grammar writer's workbench. Technical report, Xerox PARC (1996)
5. Meurers, W.D.: On implementing an hpsg theory – aspects of the logical architecture, the formalization, and the implementation of head-driven phrase structure grammars. In: Erhard W. Hinrichs, W. Detmar Meurers, and Tsuneko Nakazawa: *Partial-VP and Split-NP Topicalization in German – An HPSG Analysis and its Implementation*. Arbeitspapiere des SFB 340 Nr. 58, Universität Tübingen (1994)
6. XTAG Research Group: A lexicalized tree adjoining grammar for english. Technical Report IRCS-01-03, IRCS, University of Pennsylvania (2001)
7. Abeillé, A.: Une grammaire d'arbres adjoints pour le français. Editions du CNRS, Paris (2002)
8. Koenig, J.P., Jurafsky, D.: Type underspecification and on-line type construction in the lexicon. In: Proceedings of WCCFL94. (1995)
9. Xia, F.: Automatic Grammar Generation from two Different Perspectives. PhD thesis, University of Pennsylvania (2001)
10. Candito, M.H.: Organisation Modulaire et Paramétrable de Grammaires Electroniques Lexicalisées. PhD thesis, Université de Paris 7 (1999)
11. Perrier, G.: Les grammaires d'interaction. Université Nancy 2 (2003) Habilitation à diriger des recherches.
12. Kahane, S.: Grammaires d'unification polarisées. In: Proc. TALN 2004, Fès (2004)
13. Perlmutter, D.: Surface structure constraints in syntax. *Linguistic Inquiry* **1** (1970) 187–255
14. Kroch, A., Joshi, A.K.: The linguistic relevance of tree adjoining grammar. Technical report, IRCS, Philadelphia (1985)
15. Duchier, D., Niehren, J.: Dominance constraints with set operators. In: Proceedings of the First International Conference on Computational Logic (CL2000). Volume 1861 of Lecture Notes in Computer Science., Springer (2000) 326–341