

XMG: A Modular MetaGrammar Compiler

Simon Petitjean

Univ. Orleans, LIFO EA 4022, F-45067 Orleans, France

Abstract. XMG (eXtensible MetaGrammar) is a metagrammar compiler which has already been used for the design of large scale Tree Adjoining Grammars and Interaction Grammars. Due to the heterogeneity in the field of grammar development (different grammar formalisms, different languages, etc), a particularly interesting aspect to explore is modularity. In this paper, we discuss the different spots where this modularity can be considered in a grammar development, and its integration to XMG.

1 Introduction

1.1 Grammar Engineering

Nowadays, a lot of applications have to deal with languages and consequently need to manipulate their descriptions. Linguists are also interested in these kinds of resources, for study or comparison. For these purposes, formal grammars production has become a necessity. Our work focuses on large scale grammars, that is to say grammars which represent a significant part of the language.

The main issue with these resources is their size (thousands of structures), which causes their production and maintenance to be really complex and time consuming tasks. Moreover, these resources have some specificities (language, grammatical framework) that make each one unique.

Since a handwriting of thousands of structures represents a huge amount of work, part of the process has to be automatized. A totally automatic solution could consist in an acquisition from treebanks, which is a widely used technique. Semi automatic approaches are alternatives that give an important role to the linguist: they consist in building automatically the whole grammar from information on its structure. The approach we chose is based on a description language, called metagrammar [1]. The idea behind metagrammars is to capture linguistic generalization, and to use abstractions to describe the grammar.

1.2 Metagrammars for Tree Adjoining Grammars

The context that initially inspired metagrammars was the one of Tree Adjoining Grammars (TAG) [2]. This formalism consists in tree rewriting, with two specific rewriting operations: adjunction and substitution. An adjunction is the replacement of an internal node by an auxiliary tree (one of its leaf nodes is labelled with \star and called foot node) with root and foot node having the same syntactic category as the internal node. A substitution is the replacement of a leaf node

(marked with \downarrow) by a tree with a root having the same syntactic category as this leaf node.

TAG is said to have an extended domain of locality, because the adjunction operation and the depth of the trees allow to represent long distance relations between nodes: two nodes of the same elementary tree can after derivation end up at an arbitrary distance from each other. Here, we will only manipulate LTAG (lexicalized-TAG), which means each elementary tree is associated with at least one lexical element.

LTAG is traditionnaly used with respect to the Condition on Elementary Tree Minimality from [3], which means that an elementary tree only encapsulates the arguments of its anchor, recursion being factored away.

What can we do to lower the amount of work implied by the conception of the grammar ? Let us take a look at some rules:

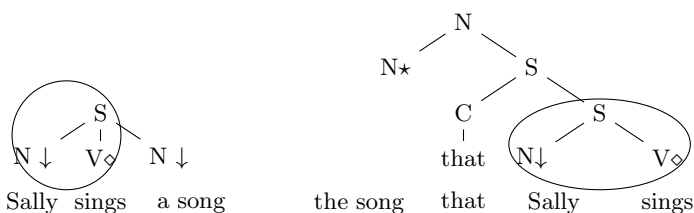


Fig. 1. Verb with canonical subject and canonical or extracted object

Those two trees share some common points: part of the structure is the same (the subject is placed before the verb in both circled parts), and the agreement constraints, given in feature structures associated to nodes (not represented here), are similar. This kind of redundancy is one of the key motivations for the use of abstractions. These abstractions are descriptions of the redundant fragments we can use everywhere they are needed.

Metagrammars are based on the manipulation of those linguistic generalizations. They consist in generating the whole grammar from an abstract description, permitting to reason about language at an abstract level.

1.3 A Need for Modularity

The metagrammatical language we will deal with here is XMG (eXtensible MetaGrammar)¹, introduced in [4]. A new project, XMG-2², started in 2010 to achieve the initial goal of the compiler, extensibility, which has not been realized yet: XMG-1 only supports tree based grammars (two formalisms, Tree Adjoining

¹ <https://sourcesup.cru.fr/xmg/>

² <https://launchpad.net/xmg>

Grammars and Interaction Grammars), and includes two levels of description, the syntactic one and the semantic one.

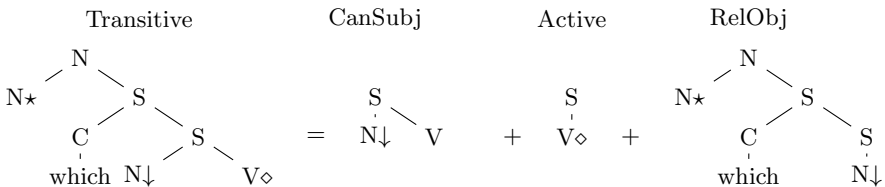
Using this metagrammatical approach for the generation of another type of linguistic resource implies the creation of a new XMG compiler. This compiler needs to provide dedicated description languages for the needed structures. A high level of flexibility is needed so that the user can assemble by their own a new metagrammatical framework.

Our goal is to go towards two levels of modularity: we want it to be possible to assemble a grammar in a modular way, thanks to a metagrammar assembled in a modular way. The first level of modularity, provided by a compiler, allows to combine abstractions to build a linguistic resource. The second one allows to build new compilers dedicated to new grammar engineering tasks.

We will begin pointing out the modularity on the grammar side in section 2. In section 3, we will focus on a new level of modularity, a metagrammatical one. In section 4, we will give an overview of what has been done, and what remains to be done. Finally, we will conclude and give some perspectives.

2 Assembling Grammars in a Modular Way

XMG consists in defining fragments of the grammar, and controlling how these fragments can combine to produce the whole grammar. The following figure shows the intuition of the combination of fragments to produce a tree for transitive verbs. It is done by combining three tree fragments, one for the subject (in its canonical form, that we noticed redundant previously), one for the object (relative) and one for the active form.



To build a lexicon, the metagrammar is first executed in a non-deterministic way to produce descriptions. Then these descriptions are solved to produce the models which will be added to the lexicon.

2.1 The Control Language and the Dimension System

The main particularity of XMG is that it allows to see the metagrammar as a logical program, using logical operators.

The abstractions (possibly with parameters) we manipulate are called classes. They contain conjunctions and disjunctions of descriptions (tree fragments descriptions for TAG), or calls to other classes. This is formalized by the following control language:

$$\begin{aligned}
 \textit{Class} & := \textit{Name}[p_1, \dots, p_n] \rightarrow \textit{Content} \\
 \textit{Content} & := \langle \textit{Dim} \rangle \{ \textit{Desc} \} \mid \textit{Name}[\dots] \mid \textit{Content} \vee \textit{Content} \\
 & \quad \mid \textit{Content} \wedge \textit{Content}
 \end{aligned}$$

For example, we can produce the two trees of the figure 1 by defining the tree fragments for canonical subject, verbal morphology, canonical object and relativized object, and these combinations:

$$\begin{aligned}
 \textit{Object} & \rightarrow \textit{CanObj} \vee \textit{RelObj} \\
 \textit{Transitive} & \rightarrow \textit{CanSubj} \wedge \textit{Active} \wedge \textit{Object}
 \end{aligned}$$

This part of metagrammar says that an object can either be a canonical object or a relative object, and that the transitive mode is created by getting together a canonical subject, an active form and one of the two object realizations.

Notice that descriptions are accumulated within dimensions, which allow to separate types of data. Sharing is still possible between dimensions, by means of another dimension we call interface. In XMG's TAG compiler for example, the *syn* dimension accumulates tree descriptions while the *sem* dimension accumulates predicates representing the semantics. Each dimension comes with a description language, adapted to the type of data it will contain. For each type of description we need to accumulate, we have to use a different description languages. The first version of XMG provides a tree description language (for TAG or Interaction Grammars) associated with the *syn* dimension and a language for semantics associated with the *sem* dimension.

A Tree Description Language. For trees in TAG, we use the following tree description language:

$$\begin{aligned}
 \textit{Desc} & := x \rightarrow y \mid x \rightarrow^+ y \mid x \rightarrow^* y \mid x \prec y \mid x \prec^+ y \mid x \prec^* y \mid x[f:E] \\
 & \quad \mid x(p:E) \mid \textit{Desc} \wedge \textit{Desc}
 \end{aligned}$$

where x and y are node variables, \rightarrow and \prec dominance and precedence between nodes ($+$ and $*$ respectively standing for transitive and reflexive transitive closures). $'.'$ is the association between a property p or a feature f and an expression E . Properties are constraints specific to the formalism (the fact that a node is a substitution node for example), while features contain linguistic information, such as syntactic categories, number or gender.

When accumulated, the tree description in the syntactic dimension is still partial. The TAG elementary trees that compose the grammar are the models for this partial description. They are built by a tree description solver, based on constraints to ensure the well-formedness of the solutions. XMG computes

minimal models, that is to say models where only the nodes of the description exist (no additional node is created).

Here is a toy metagrammar, composed of three description classes (representing canonical subject, relative object, active form) and one combination class (transitive mode):

$$\begin{aligned}
\text{CanSubj} &\rightarrow \langle \text{syn} \rangle \{ (s_1[\text{cat} : S] \rightarrow v_1[\text{cat} : V]) \wedge (s_1 \rightarrow n_1(\text{mark} : \text{subst})[\text{cat} : N]) \\
&\quad \wedge (n_1 \prec v_1) \} \\
\text{RelObj} &\rightarrow \langle \text{syn} \rangle \{ (n_2[\text{cat} = N] \rightarrow n_3(\text{mark} = \text{adj})[\text{cat} = N]) \wedge (n_2 \rightarrow s_2[\text{cat} = S]) \\
&\quad \wedge (n_3 \prec s_2) \wedge (s_2 \rightarrow c) \wedge (s_2 \rightarrow s_1[\text{cat} = S]) \wedge (c \prec s_1) \\
&\quad \wedge (c \rightarrow \text{wh}[\text{cat} = \text{wh}]) \wedge (s_1 \rightarrow n_1[\text{cat} = n]) \} \\
\text{Active} &\rightarrow \langle \text{syn} \rangle \{ (s_1 \rightarrow v_2[\text{cat} : V]) \} \\
\text{Transitive} &\rightarrow \text{CanSubj} \wedge \text{RelObj} \wedge \text{Active}
\end{aligned}$$

The minimal models for the classes named CanSubj, Active and Object are the trees with matching names on the previous figure. The tree Transitive is a minimal model for the description accumulated in class Transitive.

A Language for Semantics. To describe semantics, we use another description language, which is:

$$\text{SemDesc} := \ell : p(E_1, \dots, E_n) \mid \neg\ell : p(E_1, \dots, E_n) \mid E_i \ll E_j \mid E$$

where ℓ is a label for predicate p (of arity n) and \ll is a scope-over relation for dealing with quantifiers. To add binary relations to the semantic dimension, we can use a class of this type:

$$\text{BinaryRel}[Pred, X, Y] \rightarrow \langle \text{sem} \rangle \{ Pred(X, Y) \}$$

When instantiated with $Pred = \text{love}$, $X = \text{John}$, $Y = \text{Mary}$, calling the class *BinaryRel* accumulates the predicate $\text{love}(\text{John}, \text{Mary})$.

2.2 Principles

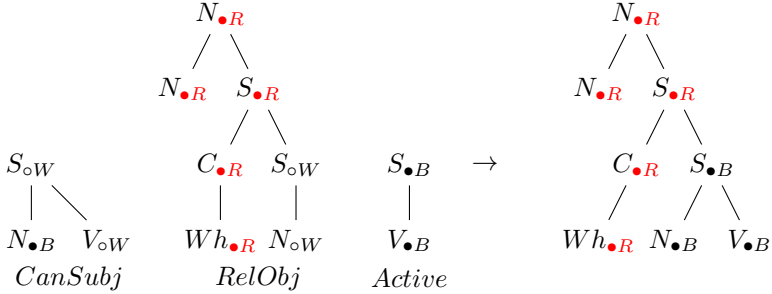
Some additional sets of constraints we call principles are available. Their goal is to check some properties in the resulting models of the compilation, they are consequently dependent from the target formalism. For example, in TAG, the color principle is a way to forbid some fragments combination, by associating colors to each node.

A valid model is a model in which every node is colored either in red or black. When unifying nodes, their colors are merged: a red node must not unify, a white node has to unify with a black node, creating a black node, and a black node can only unify with white nodes. The only valid models are the ones in which every node is colored either in red or black. The following table shows the results of colors unifications.

For example, if we consider our previous example, the colored trees of the metagrammar are the following:

| | | | | |
|----------------|----------------|----------------|----------------|---|
| | ● _B | ● _R | ○ _W | ⊥ |
| ● _B | ⊥ | ⊥ | ● _B | ⊥ |
| ● _R | ⊥ | ⊥ | ⊥ | ⊥ |
| ○ _W | ● _B | ⊥ | ○ _W | ⊥ |
| ⊥ | ⊥ | ⊥ | ⊥ | ⊥ |

Fig. 2. Unification rules for colors



The tree description solver (ignoring the colors) will produce models where the nodes labelled *S* of *CanSubj* and *Active* unify with any of the two nodes labelled *S* in *RelObj*, where the nodes labelled *V* do not unify, etc. But when filtering with the colors principle, the only remaining model is the one of the right, which is linguistically valid, contrary to the others.

We can also cite the rank principle: we use it to add constraints on the ordering of nodes in the models of the description. In French for example, clitics are necessarily ordered, so we associate a rank property to some nodes, with values that will force the right order.

3 Assembling Metagrammars in a Modular Way

The main aim of the XMG-2 project is to make it possible for the linguist to design new metagrammatical scopes, that can accommodate a large number of linguistic theories. A modular way to realize this ambition is to provide a set of bricks the user can pick or create and combine to build the compiler he needs. Those bricks could be used to design new description languages, new principles, etc.

3.1 A Modular Architecture

XMG compiler comes with a modular processing chain. This chain is composed of two phases. The first one consists in translating the metagrammatical description into executable code.

Tokenizer → Parser → Type Checker → Unfolder → Code Generator

Fig. 3. Compilation steps

First, the description is analysed and turned into an abstract syntax tree. The types into this tree are checked. The tree is then unfolded into terms of depth one, representing instructions. Instructions are finally translated into code.

The second phase corresponds to the generation of the resource.

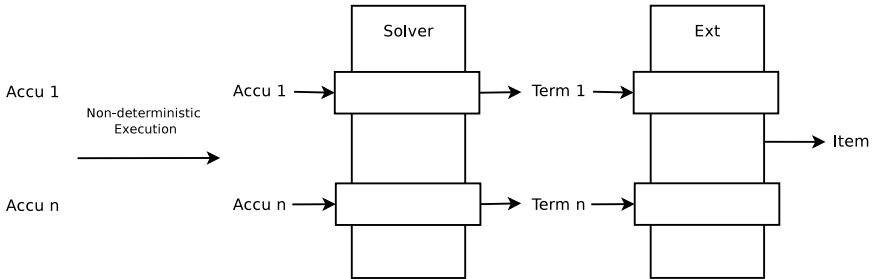


Fig. 4. Generation steps

The execution of the non-deterministic code generated by the compiler triggers accumulations in the dimensions. Each accumulation is composed of structures, and of a set of constraints over these structures. A solver extracts the models from the accumulations. The terms resulting from the solving are then translated into an output language.

The particularity of XMG is to make it possible to choose the modules that suits the best the user's metagrammar. By this mean, descriptions accumulated in different dimensions can be handled differently. For example, the end of the processing chain for TAG is a tree description solver, that builds the grammar's elementary trees from the descriptions accumulated in the syntactic dimension. The user can choose the kind of output the compiler will produce: he can interactively observe the grammar he produced, or produce an XML description of the grammar. This description can be used by a parser (for example TuLiPA [5]³ for TAG, or LeoPar⁴ for IG).

The modules of the processing chains are contributed by the XMG-2 bricks. The new compiler includes bricks that recreate the two processing chains (for Tree Adjoining Grammars and Interaction Grammars) featured by XMG-1.

³ <https://sourcesup.cru.fr/tulipa/>

⁴ <http://wikilligramme.loria.fr/doku.php?id=leopar:leopar>

3.2 Representation Modules

As we wish to build a tool which is as universal as possible, being independent from the formalism is a priority. To achieve this goal, we need to be able to describe a large number of types of structure into XMG. We saw the dimension system was useful to separate syntax from semantics, but adding new dimensions also allows to describe and combine other levels of description. A set of dimensions, with description language, has recently been proposed.

These dimensions are packaged into XMG-2 bricks and can be used to build new compilers. Different dimensions can be built from similar sets of bricks: for example, feature structures, which can be used in a lot of formalisms, are provided by a brick. Getting the support for feature structures inside a new dimension can be done simply by plugging the feature structure brick into the new dimension brick.

Syntactic Dimensions. In [6], description languages for two syntactic formalisms, namely Lexical Functional Grammars (LFG) and Property Grammars (PG), are proposed. Here, we will focus on Property Grammars, because they differ from TAG in many aspects. PG are not based on tree rewriting but on a local constraints system: the properties. A property concerns a node and applies constraints over its children nodes. One of the interesting aspects of PG is the ability to analyse ungrammatical utterances. When parsing a utterance, its grammaticality score is lowered at every violated property. Here, we will consider these six properties:

| | | |
|---------------------|-----------------------|---|
| Obligation | A: ΔB | at least one B child |
| Uniqueness | A: B! | at most one B child |
| Linearity | A: B<C | B child precedes C child |
| Requirement | A: B \Rightarrow C | if a B child, then also a C child |
| Exclusion | A: B \nRightarrow C | B and C children are mutually exclusive |
| Constituency | A: S | children must have categories in S |

A real size PG consists in an inheritance hierarchy of linguistic constructions. These constructions are composed of feature structures and a set of properties. Variables are manipulated on both sides, and can be used to share data between them. Figure 5 represents a part of the hierarchy built in [7] for French.

The V-n construction of the figure says that in verbs with negation in French, negation implies the presence of an adverb *ne* labelled with category *Adv - ng* (*ne*) and/or an adverb labelled with category *Adv - np* (like *pas*). We also have a uniqueness obligation over these adverbs, and an linear order must be respected (*ne* must come before *pas*). When the mode of the verb is infinitive, the verb must be placed after the adverbs.

| | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
|--|--|--|------|---------------------------|---|--|---|--|------|--------------------------------------|------------|-------|-------------|------------------------|-----------|------------------|---|-------------------------------------|--|------|--------------------------------------|------------|--------------------|-------------|-------------------------|-----------|--|--|--------------------------------------|
| <table border="1" style="border-collapse: collapse; width: 100%;"> <tr> <td colspan="2" style="text-align: center;">V (Verb)</td> </tr> <tr> <td style="width: 50%;">INTR</td> <td>ID—NATURE [SCAT [1].SCAT]</td> </tr> <tr> <td colspan="2">const. V: [1] [CAT V SCAT \neg (aux-etre \vee aux-avoir)]</td> </tr> <tr> <td colspan="2" style="text-align: center;">V-m (Verb with modality) inherits V ; V-n</td> </tr> <tr> <td style="width: 50%;">INTR</td> <td>[SYN [INTRO [RECT [1] DEP Prep]]]</td> </tr> <tr> <td>uniqueness</td> <td>Prep!</td> </tr> <tr> <td>requirement</td> <td>[1] \RightarrowPrep</td> </tr> <tr> <td>linearity</td> <td>[1] \precPrep</td> </tr> </table> | V (Verb) | | INTR | ID—NATURE [SCAT [1].SCAT] | const. V: [1] [CAT V SCAT \neg (aux-etre \vee aux-avoir)] | | V-m (Verb with modality) inherits V ; V-n | | INTR | [SYN [INTRO [RECT [1] DEP Prep]]] | uniqueness | Prep! | requirement | [1] \Rightarrow Prep | linearity | [1] \prec Prep | <table border="1" style="border-collapse: collapse; width: 100%;"> <tr> <td colspan="2" style="text-align: center;">V-n (Verb with negation) inherits V</td> </tr> <tr> <td style="width: 50%;">INTR</td> <td>[SYN [NEGA [RECT [1] DEP Adv-n]]]</td> </tr> <tr> <td>uniqueness</td> <td>Adv-ng! Adv-np!</td> </tr> <tr> <td>requirement</td> <td>[1] \RightarrowAdv-n</td> </tr> <tr> <td>linearity</td> <td>Adv-ng \prec [1] Adv-ng \prec Adv-np Adv-np \prec [1].[MODE inf]</td> </tr> <tr> <td></td> <td>[1].[MODE \neginf] \prec Adv-np</td> </tr> </table> | V-n (Verb with negation) inherits V | | INTR | [SYN [NEGA [RECT [1] DEP Adv-n]]] | uniqueness | Adv-ng! Adv-np! | requirement | [1] \Rightarrow Adv-n | linearity | Adv-ng \prec [1] Adv-ng \prec Adv-np Adv-np \prec [1].[MODE inf] | | [1].[MODE \neg inf] \prec Adv-np |
| V (Verb) | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| INTR | ID—NATURE [SCAT [1].SCAT] | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| const. V: [1] [CAT V SCAT \neg (aux-etre \vee aux-avoir)] | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| V-m (Verb with modality) inherits V ; V-n | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| INTR | [SYN [INTRO [RECT [1] DEP Prep]]] | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| uniqueness | Prep! | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| requirement | [1] \Rightarrow Prep | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| linearity | [1] \prec Prep | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| V-n (Verb with negation) inherits V | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| INTR | [SYN [NEGA [RECT [1] DEP Adv-n]]] | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| uniqueness | Adv-ng! Adv-np! | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| requirement | [1] \Rightarrow Adv-n | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| linearity | Adv-ng \prec [1] Adv-ng \prec Adv-np Adv-np \prec [1].[MODE inf] | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| | [1].[MODE \neg inf] \prec Adv-np | | | | | | | | | | | | | | | | | | | | | | | | | | | | |

Fig. 5. Fragment of a PG for French (basic verbal constructions)

To describe a PG, we need to be able to represent encapsulations, variables, feature structures, and properties. We can notice that XMG classes can be seen as encapsulations, and that variables and feature structures were already used for TAG descriptions. Considering that, the XMG description language for PG can be formalized this way:

$$\begin{aligned}
 Desc_{PG} &:= x = y \mid x \neq y \mid [f:E] \mid \{P\} \mid Desc_{PG} \wedge Desc_{PG} \\
 P &:= A : \Delta B \mid A : B! \mid A : B \prec C \mid A : B \Rightarrow C \mid A : B \not\Rightarrow C \mid A : B
 \end{aligned}$$

where x, y correspond to unification variables, $=$ to unification, \neq to unification failure, $:$ to association between the feature f and some (possibly complex) expression E , and $\{P\}$ to a set of properties. Note that E and P may share unification variables.

The translation of the linguistic construction for V-m in XMG would be:

$$\begin{aligned}
 V-m &\rightarrow (Vclass \vee V-n) \wedge \langle PG \rangle \{ [INTR:[SYN:[INTRO:[RECT:X, DEP:Prep]]] \} \\
 &\quad \wedge (V : Prep!) \wedge (V : X \Rightarrow Prep) \wedge (V : X \prec Prep)
 \end{aligned}$$

Here, inheritance is made possible by calls of classes. The control language even allows to do disjunctive inheritance, like it happens in class V-m. The end of the compilation process for PG will differ from TAG's one. We don't need any solver for descriptions, the accumulation into PG dimension is the grammar. To get the properties solved for a given sentence, the solution is to use a parser as a post processor for the compiler.

Morphological Dimension. For the needs of the study of verbal morphology in Ikota [8], a morphological dimension based on the notion of topological fields [9] was proposed. The description language available inside this dimension is the following:

$$Desc_{Morph} := f \leftarrow c \mid attr = val \mid Desc_{Morph} \wedge Desc_{Morph}$$

where f is a field, declared for the whole metagrammar, c is a contribution, and \leftarrow corresponds to the accumulation of a contribution into a field. $attr = val$

means that the feature composed of this pair will be part of the accumulated description.

The execution of the metagrammar starts with the ordering of fields. This solving has to be done only once in this dimension because in the chosen morphological theory, positions are fixed. For every solution of the execution of the classes, strings are accumulated into the fields, and morphosyntactic information into features.

The output of the compilation process is not a grammar strictly speaking, but a lexicon of fully inflected forms, basically obtained by concatenation of the fields contents.

Frame Semantics Dimension. A dimension handling a second formalism for semantics was proposed in [10]. The dedicated description language allows to describe frames, which are representations of mental concepts [11] and can be represented as feature structures. The unification of frames implies the unification of their types, which belong to a type hierarchy. This specific type unification is handled by the frame compiler brick. The description language for frames is the following:

$$Desc_{Frame} := f(t, [a_1 = f_1, \dots a_n = f_n]) \mid Desc_{Frame} \wedge Desc_{Frame}$$

where f is an optional variable labeling the frame, $a_1 \dots a_n$ are attributes of the frame, and $f_1 \dots f_n$ are frame associated to these attributes. The execution of the frame dimension leads to the accumulation and combination of frame fragments.

The output for this dimension is a set of frames, that should be associated to syntactic structures. The interface between TAG trees and frames is discussed in [12].

Including a specific representation module to the XMG-1 compiler could be seen as an ad-hoc solution. This is why allowing the linguist to build their own dimension, beginning with the choice of a description language, is a central feature of the new version of XMG. A XMG-2 brick corresponding to a new representation module is composed of the definition of the language used by the brick (the dedicated description language) and of the compilation modules to handle this language.

3.3 Specific Virtual Machines

During the generation of the linguistic resource, objects corresponding to the described structures are manipulated. The main operation between structures is unification, triggered explicitly (by using the equal sign) or implicitly (by importing variables from other classes). For most of these structures, standard unification is adequate, but for some of them, specific engines have to be used. For example, feature structures (like the ones used in TAG) need a dedicated unification algorithm, corresponding to set union.

A XMG-2 brick for a new description language has to include the set of specific virtual machines needed to handle the unification of its structures. For

the frame semantics dimension for example, a dedicated virtual machine handling the unification of typed feature structures is contributed by the brick.

3.4 Principle Bricks

The notion of principles defined in XMG was too restrictive for our aims. Their specificity for the target formalism, for example, is incompatible with the multi-formalism ambition. An interesting way to handle principles is the one of [13], both allowing the linguist to create his own principles or to use a subset of the ones already defined. An example is the tree principle, which states that the solution models must be trees.

What we aim to provide is a meta-principles library: generic and parametrizable principles the user can pick and configure. For example, the color principle provided for TAG could be an implementation of a generic polarity principle, parametrized with the table of figure 2. Another example of meta-principle is called unicity and was already implemented in XMG-1. It is used to check the uniqueness of a specific attribute-value pair in each solution, and thus is not specific to any linguistic theory.

Principles are also packaged into XMG-2 bricks. This means that for any new metagrammatical scope where trees have to be solved, the tree principle brick just has to be plugged into the new (or existing) dimension brick.

For the morphological dimension discussed early, a principle brick handling linear ordering constraints between fields was created.

3.5 Dynamic Definition of a Metagrammar Compiler

To build their own metagrammatical scope, one only has to create and configure the dimensions he needs and the properties he wants to check on them. Building a compiler consists in picking and combining independent modules, which we call compiler bricks. XMG-2 provides a compiler builder, that assembles the needed parts of the compiler according to a description of the connections between the bricks. The tokenizer and the parser for the metagrammatical language are automatically generated from this description, and each brick contributes its own compilation modules.

One of the main advantages of this modular approach is that the specific part of the compiler is mostly written automatically, and new features could be added just for experiments. A user can either use an existing compiler or assemble parts to build their own. Defining the principles would just consist in taking meta-principles out from the library and instantiate them.

Building a metagrammar compiler in this way allows to deal with a large range of linguistic theories, or even to quickly experiment while creating a new grammar formalism.

4 Current State of the Work

XMG project started in 2003 with a first tool, that has been used to produce large TAG grammars for French [14], German [15] and English, and a large

Interaction Grammar for French [16]. The compiler was written in Oz/Mozart, a language which is not maintained any more and not compatible with today's architectures (64 bits). It was also important to restart from scratch, in order to build a compiler more in adequation with its ambitions : modularity and extensibility.

Consequently, a new implementation started in 2010, in YAP (Yet Another Prolog) with bindings with Gecode for constraints solving. XMG-2 is currently the tool used for modeling the syntax and morphology of various African and Creole languages, and is compatible with the previous large metagrammars. It also includes the support for the dimensions discussed in this article.

5 Conclusion

In this paper, we showed how modularity, together with a metagrammatical approach, eases the development of a large scale grammar. This modularity is essential for reaching the main goal of XMG, that is to say extensibility. Getting to that means taking a big step towards multi-formalism and multi-language grammar development, and then offers new possibilities for sharing data between different types of grammar, or even for comparing them.

Two levels of modularity are given by XMG. The first one is the grammatical modularity, which makes it easier to generate and maintain large scale grammars thanks to the definition and the combination of abstractions. The second level of modularity is metagrammatical: XMG-2 provides a way to build new compilers by defining and combining elementary parts of compiler, called compiler bricks. The users have different options: they can use existing compilers (the one for TAG and 'flat' semantics for example), combine bricks to build a new type of compiler (like a compiler having two TAG dimensions, for two different languages), or create their own bricks, to combine them with existing ones (a brick for dependency grammars for example).

References

1. Candito, M.: A Principle-Based Hierarchical Representation of LTAGs. In: Proceedings of COLING 1996, Copenhagen, Denmark (1996)
2. Joshi, A.K., Schabes, Y.: Tree Adjoining Grammars. In: Rozenberg, G., Salomaa, A. (eds.) Handbook of Formal Languages. Springer, Berlin (1997)
3. Frank, R.: of Pennsylvania. Institute for Research in Cognitive Science, U.: Syntactic Locality and Tree Adjoining Grammar: Grammatical, Acquisition and Processing Perspectives. IRCS report. University of Pennsylvania, The Institute for Research in Cognitive Science (1992)
4. Duchier, D., Le Roux, J., Parmentier, Y.: The Metagrammar Compiler: An NLP Application with a Multi-paradigm Architecture. In: Van Roy, P. (ed.) MOZ 2004. LNCS, vol. 3389, pp. 175–187. Springer, Heidelberg (2005)

5. Kallmeyer, L., Lichte, T., Maier, W., Parmentier, Y., Dellert, J., Evang, K.: TuLiPA: Towards a Multi-Formalism Parsing Environment for Grammar Engineering. In: *Coling 2008: Proceedings of the Workshop on Grammar Engineering Across Frameworks*, Manchester, England, pp. 1–8. *Coling 2008 Organizing Committee (2008)*
6. Duchier, D., Parmentier, Y., Petitjean, S.: Cross-framework Grammar Engineering using Constraint-driven Metagrammars. In: *CSLP 2011*. Karlsruhe, Allemagne (2011)
7. Guénot, M.L.: *Éléments de grammaire du français pour une théorie descriptive et formelle de la langue*. PhD thesis, Université de Provence (2006)
8. Duchier, D., Magnana Ekoukou, B., Parmentier, Y., Petitjean, S., Schang, E.: Describing Morphologically-rich Languages using Metagrammars: a Look at Verbs in Ikota. In: *Workshop on “Language Technology for Normalisation of Less-resourced Languages”, 8th SALT MIL Workshop on Minority Languages and the 4th Workshop on African Language Technology*, Istanbul, Turkey (2012)
9. Stump, G.T.: On the theoretical status of position class restrictions on inflectional affixes. In: Booij, G., van Marle, J. (eds.) *Yearbook of Morphology 1991*, pp. 211–241. Kluwer (1992)
10. Lichte, T., Diez, A., Petitjean, S.: Coupling Trees and Frames through XMG. In: *ESSLLI 2013 Workshop on High-level Methodologies for Grammar Engineering (HMGE 2013)*, Duesseldorf, Germany (2013)
11. Fillmore, C.J.: Frame semantics. In: *The Linguistic Society of Korea. Linguistics in the Morning Calm*, pp. 111–137. Hanshin Publishing (1982)
12. Kallmeyer, L., Osswald, R.: Syntax-driven semantic frame composition in Lexicalized Tree Adjoining Grammar. *Journal of Language Modelling* 1, 267–330 (2013)
13. Debusmann, R.: *Extensible Dependency Grammar: A Modular Grammar Formalism Based On Multigraph Description*. PhD thesis, Saarland University (2006)
14. Crabbé, B.: *Représentation informatique de grammaires fortement lexicalisées: Application à la grammaire d’arbres adjoints*. PhD thesis, Université Nancy 2 (2005)
15. Kallmeyer, L., Lichte, T., Maier, W., Parmentier, Y., Dellert, J.: Developing a tt-mctag for german with an rcg-based parser. In: *LREC. ELRA* (2008)
16. Perrier, G.: *A French Interaction Grammar*. In: *RANLP, Borovets, Bulgaria* (2007)