

Using Lexicon-Grammar tables for French verbs in a large-coverage parser

Elsa Tolone¹, Benoît Sagot²

1. Institut Gaspard Monge – Université Paris-Est
5 boulevard Descartes – Champs-sur-Marne
77454 Marne-la-Vallée Cedex 2 – France
elsa.tolone@univ-paris-est.fr

2. ALPAGE, INRIA Paris-Rocquencourt & Université Paris 7
Domaine de Voluceau – Rocquencourt B.P. 105
78153 Le Chesnay Cedex – France
benoit.sagot@inria.fr

Abstract

In this paper, we describe the integration of Lexicon-Grammar tables for French verbs in the large-coverage FRMG parser and the evaluation of the resulting parser. This integration required a conversion step so as to extract the syntactic information encoded in Lexicon-Grammar tables and represent it in the NLP lexical formalism used by FRMG, i.e., the Alexina framework (that of the *Lefff* lexicon, on which the standard version of FRMG relies). We describe the linguistic basis of this conversion process, and the resulting lexicon. We compare the results of the FRMG parser on the EASy reference corpus depending on whether it relies on the verb entries of the *Lefff* or those of the converted Lexicon-Grammar verb tables.

1. Introduction

Lexicon-Grammar tables are currently one of the major sources of syntactic lexical information for the French language. Moreover, several Lexicon-Grammar tables exist for other languages, such as Italian, Portuguese, Modern Greek, Korean, and others. Their development was initiated as early as the 1970s by Maurice Gross, at the LADL and then the IGM (Université Paris-Est) (Gross, 1975; Boons et al., 1976; Guillet and Leclère, 1992). Lexical information is represented in the form of *tables*. Each table puts together elements of a given category (for a given language) that share a certain number of *defining features*, which usually concern sub-categorization. These elements form a *class*. These tables are represented as matrices: each row corresponds to a lexical item of the corresponding class; each column lists all features¹ that may be valid or not for the different members of the class; at the intersection of a row and a column, the symbol + (resp. –) indicates that the feature corresponding to the column is valid (resp. not valid) for the lexical entry corresponding to the row. As far as the French language is concerned, 61 tables for simple verbs have been developed, as well as 59 tables for predicative nouns, 65 tables for idiomatic expressions (mostly verbal), and 32 tables for (simple and idiomatic) adverbs.

Current tables suffer from various types of inconsistency and incompleteness. In particular, defining features are not represented in the tables.² To remedy this situation, *tables of classes* are being developed at IGM for each category, and notably for verbs, which associate the set of their defining features with each class (Paumier, 2003). Preliminary results of this long-term effort allowed us to convert verb tables into a format suitable for their use within a

large-scale parser for French, the FRMG parser (Thomasset and de La Clergerie, 2005). This format is that of the Alexina framework, in which the lexicon used by the standard FRMG was developed. This lexicon is the *Lefff* (see below).

This paper is organized as follows. Section 2 briefly describes the *Iglex* lexicon. Section 3 introduces the Alexina format and the *Lefff* NLP syntactic lexicon for French. Section 4 provides an overview of the interpretation and conversion process that allowed us to build an Alexina version of *Iglex*. Then section 5 describes how we coupled this converted lexicon with the FRMG parser, and section 6 compares the results of the FRMG parser on the EASy reference corpus depending on whether it relies on the verb entries of the *Lefff* or those of the converted *Iglex* (i.e., on converted Lexicon-Grammar verb tables). Finally, several further steps for this work are evoked in section 7.

2. The verbal lexicon *Iglex*

A table of classes groups the list of all syntactic features identified for the corresponding category as columns, and the rows list all classes defined for this category. At the intersection of a row and a column, the symbol + (resp. –) indicates that the corresponding feature is valid (resp. not valid) for all elements of the class (i.e., for all entries of the corresponding table). The symbol *o* indicates that the feature is explicitly coded in the corresponding table, because it is valid only for some of its entries. Finally, the symbol ? means that this cell has not been filled in yet.

The development of the table of verb classes and that of noun classes is close to completion (Constant and Tolone, 2008), since the coding ? is now only used when a given feature has not yet been studied for a given class. Thanks to this work which makes all syntactic features of verbs in the Lexicon-Grammar tables consistent and explicit, it was possible to build a structured version of the tables, available in text or XML format, and called the *Iglex* lexicon

¹For example, the feature $N_0 V$ means “possible head of an intransitive construction with initial subject noun phrase”; the feature [passif] means “passive diathesis possible”.

²This also motivated the work described in (Gardent et al., 2005). A comparison between the textual version of the tables, that is used in the present work, and the work of (Gardent et al., 2005) can be found in (Constant and Tolone, 2008).

(Constant and Tolone, 2008).³ The construction of *Iglex* relies on the *LGExtract* tool, that takes as input the tables of a given category, the corresponding table of classes and a configuration file. This file defines how each feature (as extracted from the table of classes, or, in the case of features that are coded *o*, extracted from the corresponding table) contribute to building the *Iglex* entry.

Our conversion process starts from *Iglex* verbal entries and turns them into entries in the Alexina format, i.e., the same as the format of the syntactic lexicon *Lefff*.

3. The *Lefff* syntactic lexicon and the Alexina format

The *Lefff* (Lexique des formes fléchies du français — *Lexicon of French inflected form*) is a large-coverage syntactic lexicon for French (Sagot et al., 2006; Sagot and Danlos, 2007)⁴. It relies on the Alexina framework for the acquisition and modeling of morphological and syntactic lexicons. To represent lexical information, an Alexina lexicon relies on a two-level architecture:

- the *intensional* lexicon associates (among others) an inflection table and a canonical sub-categorization frame with each entry and lists all possible redistributions from this frame;
- the *compilation* of the intensional lexicon into an *extensional lexicon* builds different entries for each inflected form of the lemma and every possible redistribution.

For example, consider the following (simplified) intensional entry:

```
clarifier1 Lemma:v;<Suj:cln|scompl|sinf|sn,
Obj:(cl|scompl|sn)>;
%ppp_employé_comme_adj,%actif,%passif,
%se_moyen_impersonnel,%passif_impersonnel
```

It describes an entry of the verbal lemma *clarifier* (*clarify*) which is transitive (two arguments canonically realized by the *syntactic functions* *Suj* and *Obj* listed between brackets), and which allows for the functional redistributions *past participle used as an adjective*, *active* (the default distribution), *impersonal middle-voice “se” construction*, *impersonal passive*, and *passive*.

The different syntactic functions are defined in the *Lefff* by criteria close to that used by the authors of the verb valency lexicon DICOVALENCE (van den Eynde and Mertens, 2006), i.e., they rely for a large part on cliticization and other pronominal features. The *Lefff* uses the following syntactic functions: *Suj* (subject), *Obj* (direct object), *Objà* (indirect object canonically introduced by preposition “à”), *Objde* (indirect object canonically introduced by preposition “de”), *Loc* (locative), *Dloc* (delocative), *Att* (attribute), *Obl* or *Obl2* (other oblique arguments). Defining criteria for these functions are described in (Sagot and Danlos, 2007).

³Partial on-line distribution under the LGPL-LR license at <http://infolingu.univ-mlv.fr/english>, Language Resources > Lexicon-Grammar > View.

⁴On-line distribution under the LGPL-LR license at <http://gforge.inria.fr/projects/alexina/>

Each syntactic function can be realized by three types of *realizations*: *clitic pronouns*, *direct phrases* (nominal phrase (*sn*), adjectival phrase (*sa*), infinitive phrase (*sinf*), completive (*scompl*), indirect interrogative (*qcompl*)) and *prepositional phrases* (direct phrases preceded by a preposition, such as *de-sn*, *à-sinf* or *pour-sa*).⁵ Finally, a function whose realization is not optional has its realizations list between brackets.

Complementary syntactic information (control, mood for completives, etc.) are represented by *macros* (*@Ctrl-SujObj*, *@ComplSubj*, etc.) whose formal interpretation varies according to the context of use. An LFG modeling of these macros is provided with the *Lefff*.

4. Conversion of the verbal lexicon *Iglex* into a lexicon in the Alexina format

4.1. Sketch of the conversion process

Each entry in *Iglex* is associated with a set of constructions that can be classified into several types:

1. the “base” construction(s), defining feature for the originating class of the entry;
2. “extended base” constructions, obtained by adding extra arguments to the base construction; in practice, these constructions are all intermediate constructions between the base construction and a construction called “maximal extended base” construction, or MEBC;
3. constructions that are variants of the base construction, obtained by deleting one or several arguments, or by changing the realization type (e.g., Qu P can become $V^i \text{ inf } W$, as for the direct object of *savoir* — *to know* —, that can be a finite phrase but also an infinitive phrase);
4. constructions that are in fact redistributions, such as [passif de], that denotes the possibility of having a passive with an agent introduced by *de* (cf. *Pierre est aimé de Marie* — *Pierre is loved by Marie*) or $N_1 \text{ est } V_{pp} \text{ de ce Qu P}$ (cf. *Marie est étonnée de ce que Pierre soit là* — *Marie is surprised that Pierre is here*);
5. constructions that should seemingly have led to distinct entries, called “secondary entries”, such as neutral constructions of transformations like $N_1 \text{ se } V \text{ de ce Qu P}$ (cf. *Luc se félicite d’avoir réussi à séduire Léa* vs. *Max félicite Luc qu’il ait réussi à séduire Léa* — *Luc is very pleased he succeeded in seducing Léa* vs. *Max congratulates Luc for having succeeded in seducing Léa*).

We developed a method for *aligning* two constructions, i.e., for building correspondences between arguments despite their surface differences⁶ and their possible deletion.

⁵*à-scompl* and *de-scompl* represent realizations of the form *à/de ce que P*.

⁶For example, Qu P vs. N1, or à N1 vs. Prép N1 if in addition it is known that Prép can be à.

This method allows us to identify and align the MEBC and its variants, which we put together in a single entry of the final lexicon, called *canonical* entry. Among the other constructions, those that correspond to standard redistributions ([*passif par*], [*extrap*]. . .) lead to the inclusion of the corresponding redistribution in the canonical entry.⁷ Other constructions lead to the creation of extra entries, because they correspond to secondary entries (5th type in the enumeration above) or because they involve redistributions that have not yet been identified in the Alexina format.

Once the entries to be produced are identified, we build sub-categorization frames. First, we build the frame corresponding to the maximal construction for each entry (the MEBC for the canonical entry, and their unique construction for secondary entries). The syntactic function of each argument is obtained by the following heuristics. First, the first argument always receives the function *Suj* (subject). The first post-verbal argument, if it is direct, receives the function *Obj*, apart from entries of table 32NM. Then, an argument introduced by *à* (resp. *de*) receives the syntactic function *Objà* (resp. *Objde*), except if an additional indicator contradicts this choice (e.g., for an N_1 argument introduced by *à*, the feature $\grave{a} N_1 = Ppv =: le$ shows it must receive the syntactic function *Obj*, as in *Il apprend à conduire / Il l'apprend — He is learning how to drive / He is learning it*). Arguments introduced by *LOC* have the syntactic function *Loc*, except those of the form $LOC N_i$ source or for which $LOC N_i =: de N_i$ source is a valid feature, which receive the syntactic function *Dloc*. Finally, other arguments are considered as *Att* if they are direct, and as *Obl* if they are introduced by a preposition (*Obl2* if an *Obl* already exists).

The realizations of these syntactic functions are built in two steps. First, the kind of phrase (nominal, infinitive, etc.) is determined. Then, possible introducers are extracted from the set of corresponding prepositions and other introducers (e.g., *et — and*). For the canonical entry, all variants of the MEBC lead to modifications of the resulting sub-categorization frame, by adding realizations and making some arguments optional.

Other types of information are then added so as to finalize the entry, such as the originating table and the corresponding row number, as well as a frequency information extracted from the DELA. Finally, syntactic macros concerning the auxiliary, the mood of completive arguments, idiomatic clitics (*se, en, ne, etc.*) and control phenomena are extracted and added to the final entry.

⁷In the table of classes, the feature [*passif par*] (the standard passivability) is not yet correctly described, even for transitive classes. Considering this lack of information as a negative information (non-passivable), as done for other features, leads to a really incomplete lexicon. Therefore, we decided to add the corresponding *%passif* redistribution to all entries that have an argument whose syntactic function is *Obj* (direct object). Note that direct complements of the entries of table 32NM do not receive the function *Obj* (see below). Therefore, our heuristics is valid, apart from rare cases of non-passivability such as *regarder* (often *to look at*, but also *to concern*) in the sense of *concerner* (*to concern*).

4.2. Resulting lexicon

The resulting verbal lexicon contains 16,903 entries for 5,694 distinct verb lemmas (on average, 2.96 entries per lemma). As a comparison, the *Lefff* only contains 7,072 verbal entries for 6,818 distinct verb lemmas (on average, 1.04 entries per lemma). The resulting lexicon extracted from *Iglex*, despite the fact that it describes fewer verbal lemmas, has a larger coverage in terms of syntactic constructions and therefore is much more ambiguous. At the extensional level, the *Lefff* has 361,268 entries whereas the lexicon extracted from *Iglex* has 763,555 entries.

The construction of this lexicon from *Iglex* according to the process described in this section is achieved by a *perl* script that contains less than 1,000 lines. The conversion in itself, i.e., the execution of the script of the whole *Iglex*, takes less than a minute.⁸ Therefore, if a new version of the Lexicon-Grammar French verb tables or of the corresponding table of classes is released, building the new corresponding Alexina-format lexicon is a matter of seconds, and does not require any new development.

5. Integration in the FRMG parser

The main goal of this work is to allow the use of the linguistic data coded in Lexicon-Grammar tables for French to be used as a lexical database for a French parser. Among the various parsers that rely on a syntactic lexicon in the Alexina format, we chose the FRMG parser (Thomasset and de La Clergerie, 2005). It relies on a compact factorized Tree Adjoining Grammar (TAG) generated from a meta-grammar, and on the *Lefff*. The compilation and execution of the parser is performed by the DIALOG system (de La Clergerie, 2005). The result of the parsing itself is a derivation shared forest, that undergoes a symbolic (weight-based) disambiguation process so as to output only one parse. In case a sentence is not covered by the grammar and the lexicon, FRMG outputs several partial parses that cover disjoint parts of the sentence (however, no attempt is made to reassemble these partial parses into a global parse). FRMG takes as its input the result of the presyntactic processing chain SXPipe (Sagot and Boullier, 2008), which converts a raw text into a lattice of forms known by the lexicon (namely, the *Lefff*).⁹

The integration of the Alexina-format lexicon extracted from *Iglex* in the FRMG parser is straightforward: in its standard configuration FRMG's lexer performs calls to a lexical database built from the *Lefff*. We shall call this standard parser $FRMG_{Lefff}$. What is required to use the lexical information from Lexicon-Grammar verb tables is to replace verbal entries in the *Lefff* by those of the lexicon built from *Iglex* while keeping other *Lefff* entries, to build the corresponding lexical database, and to tell FRMG to use it rather than the default *Lefff*-only one.

However, several verbal entries which are not covered by *Iglex* had to be added as well: entries for auxiliaries and

⁸On a 2.4 GHz machine using Ubuntu Linux.

⁹SXPipe includes, among others, modules for (deterministic) sentence splitting and tokenization, as well as non-deterministic spelling error correction, named entity detection and identification of compound forms.

semi-auxiliaries, some raising verbs, impersonal verbs and light verbs. The result is a variant of the FRMG parser, that we shall call FRMG_{Igllex} , to distinguish it from the standard FRMG_{Lefff} .

6. Evaluation and discussion

We evaluated both FRMG_{Lefff} and FRMG_{Igllex} by parsing the manually annotated part of the EASy corpus (Paroubek et al., 2005), i.e., 4,306 sentences of diverse genres (journalistic, medical, oral, questions, literature, and others).

We used the metrics defined and used during the first French parsing evaluation campaign EASy, which took place at the end of 2005 (Paroubek et al., 2006). These metrics rely on both (non-recursive) « chunks » and « relations » (dependencies between full words), for which the standard measures (precision, recall, f-measure) are applied. In this paper, we simply provide f-measures.

Before discussing the results of these experiments, some precautions must be taken:

- the conversion process described in this paper and its still preliminary implementation certainly contain errors, and we evaluate a variant of FRMG that relies on *converted* entries extracted from Lexicon-Grammar tables, not directly on Lexicon-Grammar entries from the tables;
- the $Lefff$ was developed in parallel with EASy campaigns, unlike Lexicon-Grammar tables; some choices in the EASy annotation guide may have influenced choices made during the development of the $Lefff$, whereas it is obviously not the case for Lexicon-Grammar tables;
- as mentioned in the previous section, $Igllex$ had to be completed by various lexical entries from the $Lefff$, but other entries may still need to be added.

Comparative results for both parsers are shown on Table 1, with detailed results for some illustrative sub-corpora. As can be seen, results are for now a bit better for FRMG_{Lefff} . We do not think that this result questions the relevance of using Lexicon-Grammar tables in a parser, especially given the above-mentioned precautions. In particular, we remain convinced that using as rich a lexical resource as possible is an efficient way to improve the quality of a parser, as has been shown for example by the results of the work described in (Sagot and de La Clergerie, 2006).

However, parsing times are more than twice as high with FRMG_{Igllex} as with FRMG_{Lefff} (median average time per sentence: 0.62 s vs. 0.26 s), which is certainly a consequence of the higher average number of entries per lemma, which is three times higher in the lexicon extracted from $Igllex$ than in the $Lefff$ (see above). In fact, these higher parsing times necessarily lead to a higher ratio of parsing failures because of reaching the timeout, which leads to the construction of partial parses whose quality can only be lower.

Nevertheless, on several sub-corpora, FRMG_{Igllex} performs better in terms of chunk f-measure; but results on

relations are better with FRMG_{Lefff} , apart from two sub-corpora. A careful study of the results shows the following interesting facts:

- FRMG_{Igllex} performs better than FRMG_{Lefff} on several relations, such as “adjective modifier” and “adverb modifier”, and also on two relations for which results are anyway quite low (“preposition modifier” and “apposition”);
- the relation “(subject of object) attribute” is that for which the difference in terms of recall is the highest (34.0% vs. 58.5%);
- the high number of verb arguments listed in $Igllex$ ’s sub-categorization frames tends to fool the usual disambiguation heuristics according to which “arguments are preferred to modifiers”: any phrase that can be parsed as a verbal argument tends to be done in this way. For example, in a sentence such as [...] *on estime que cette décision [ferait] dérailler le processus de paix* (it is estimated that this decision would derail the peace process), FRMG_{Igllex} considers *de paix* (*peace*_{genitive}) as an argument of *estimer* (*estimer qqch de qqch/qqn* — to estimate something about something/somebody), whereas FRMG_{Lefff} gets the correct parse.

In the short term, the following statement can be made. Many sentences get a full parse from FRMG_{Igllex} but not from FRMG_{Lefff} , and vice versa. For example, on the *general_lemonde* sub-corpus, 177 sentences are fully parsed by both parsers, 85 only by FRMG_{Lefff} , 76 only by FRMG_{Igllex} , and 111 by neither of them. Since experience shows that partial parses lead to worse results (approx. 10 points lower in terms of f-measure on EASy relations), an interesting experiment would be to couple both parsers in such a way that if only one of them builds a full parse for a given sentence, this parse is kept (what should be done in other cases remains an open question). The results of such a “meta-parser” should be better than those of both parsers.

In the long term, it is important to benefit from this complementarity between both resources. It will be interesting to study the differences between errors made by both parsers, in particular thanks to techniques such as those described in (Sagot and de La Clergerie, 2006). This could lead to an improvement for both resources, and in particular the lexicon converted from $Igllex$. Perhaps we will realize that most errors come from the conversion process; but some errors may come from errors in Lexicon-Grammar tables, and may therefore allow us to improve them.

7. Conclusion and future work

In this paper, we introduced a methodology and a tool for converting the textual version of Lexicon-Grammar tables into an NLP lexicon based on the Alexina framework, i.e., in the same format as the $Lefff$ syntactic lexicon for French, which is used by the FRMG parser. The relevance of the resulting lexicon is confirmed by its use for parsing the evaluation corpus of the French parsing evaluation campaign EASy.

Sub-corpus	Chunks		Relations	
	FRMG _{Lefff}	FRMG _{lglex}	FRMG _{Lefff}	FRMG _{lglex}
general_lemonde	86.8%	82.8%	59.8%	56.9%
general_senat	82.7%	83.1%	56.7%	54.9%
litteraire_2	84.7%	81.5%	59.2%	56.3%
medical_2	85.4%	89.2%	62.4%	58.6%
oral_delic_8	74.1%	73.6%	47.2%	48.5%
questions_amaryllis	90.5%	90.6%	65.6%	63.2%
<i>EASy corpus overall</i>	84.4%	82.3%	59.9%	56.6%

Table 1: EASy results for FRMG_{Lefff} and FRMG_{lglex}, expressed in terms of f-measure. For reasons of space, figures are given for the whole EASy corpus and for only a sample of sub-corpora.

The first step described here has allowed us to identify several problems in the input data (tables and tables of classes), but also several simplifications and approximations in the conversion process. Therefore, there is space for significant improvements, that could eventually lead to the construction of a syntactic lexicon for French based on Lexicon-Grammar tables. Such a lexicon would improve the quality of existing tools and resources, e.g., by fusion with other lexical resources and by integration in a large-coverage parser.

However, as mentioned in the introduction, we intend to enlarge the scope of our approach by applying the same approach on French to Lexicon-Grammar tables for other categories, but also on tables for other languages, as soon as the corresponding tables of classes become available. The next step, which should be taken soon, will deal with French predicative nouns, verbal idiomatic expressions and adverbs.

8. References

- Boons, Jean-Paul, Alain Guillet, and Christian Leclère, 1976. *La structure des phrases simples en français : Constructions intransitives*. Geneva, Switzerland: Droz.
- Constant, Matthieu and Elsa Tolone, 2008. A generic tool to generate a lexicon for NLP from Lexicon-Grammar tables. In *Proceedings of the 27th Lexis and Grammar Conference*. L'Aquila, Italy.
- de La Clergerie, Éric, 2005. DyALog: a tabular logic programming based environment for NLP. In *Proceedings of 2nd International Workshop on Constraint Solving and Language Processing (CSLP'05)*. Barcelona, Spain.
- Gardent, Claire, Bruno Guillaume, Guy Perrier, and Ingrid Falk, 2005. Maurice Gross' Grammar Lexicon and Natural Language Processing. In *Proceedings of the 2nd Language and Technology Conference (LTC'05)*. Poznań, Poland.
- Gross, Maurice, 1975. *Méthodes en syntaxe : Régimes des constructions complétives*. Paris, France: Hermann.
- Guillet, Alain and Christian Leclère, 1992. *La structure des phrases simples en français : Les constructions transitives locatives*. Geneva, Switzerland: Droz.
- Paroubek, Patrick, Louis-Gabriel Pouillot, Isabelle Robba, and Anne Vilnat, 2005. EASy : campagne d'évaluation des analyseurs syntaxiques. In *Proceedings of the EASy workshop of TALN'05*. Dourdan, France.
- Paroubek, Patrick, Isabelle Robba, Anne Vilnat, and Christelle Ayache, 2006. Data, Annotations and Measures in EASy, the Evaluation Campaign for Parsers of French. In *Proceedings of the 5th Language Resources and Evaluation Conference (LREC'06)*. Genoa, Italy.
- Paumier, Sébastien, 2003. *De la reconnaissance de formes linguistiques à l'analyse syntaxique*. Ph.D. thesis, Université Paris-Est Marne-la-Vallée, France.
- Sagot, Benoît, Lionel Clément, Éric de La Clergerie, and Pierre Boullier, 2006. The Lefff 2 syntactic lexicon for French: architecture, acquisition, use. In *Proceedings of the 5th Language Resources and Evaluation Conference (LREC'06)*. Genoa, Italy.
- Sagot, Benoît and Laurence Danlos, 2007. Améliorer un lexique syntaxique à l'aide des tables du Lexique-Grammaire : Constructions impersonnelles. *Cahiers du Cental*.
- Sagot, Benoît and Éric de La Clergerie, 2006. Error mining in parsing results. In *Proceedings of ACL/COLING'06*. Sydney, Australia: Association for Computational Linguistics.
- Sagot, Benoît and Pierre Boullier, 2008. SxPipe 2 : architecture pour le traitement présyntaxique de corpus bruts. *Traitement Automatique des Langues (T.A.L.)*, 49(2). To appear.
- Thomasset, François and Éric de La Clergerie, 2005. Comment obtenir plus des méta-grammaires. In *Proceedings of TALN'05*. Dourdan, France.
- van den Eynde, Karel and Piet Mertens, 2006. Le dictionnaire de valence DICOVALENCE : manuel d'utilisation. http://bach.arts.kuleuven.be/dicovalence/manuel_061117.pdf.