



HAL
open science

Enriching a Lexical Resource for French Verbs with Aspectual Information

Anna Kupść, Pauline Haas, Rafael Marín, Antonio Balvet

► **To cite this version:**

Anna Kupść, Pauline Haas, Rafael Marín, Antonio Balvet. Enriching a Lexical Resource for French Verbs with Aspectual Information. 3rd Conference on Language, Data and Knowledge (LDK 2021), Sep 2021, Zaragoza, Spain. 10.4230/OASICS.LDK.2021.10 . hal-03503208

HAL Id: hal-03503208

<https://hal.science/hal-03503208v1>

Submitted on 27 Dec 2021

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Enriching a Lexical Resource for French Verbs with Aspectual Information

Anna Kupść  

CLLE and Université Bordeaux Montaigne, France

Pauline Haas  

UMR Lattice 8094 and Université Paris 13, France

Rafael Marín 

UMR STL 8163, CNRS and Université de Lille, F-59000, France

Antonio Balvet  

UMR STL 8163, CNRS and Université de Lille, F-59000, France

Abstract

The paper presents a syntactico-semantic lexicon of over a thousand French verbs. It has been created by manually adding lexical aspect features to verb frames from TreeLex [16]. We present how the original syntactic resource has been adapted to the current project, our aspect assignment procedure and an overview of the resulting lexical resource.

2012 ACM Subject Classification Computing methodologies → Language resources; Computing methodologies → Lexical semantics; Computing methodologies → Information extraction

Keywords and phrases computational semantics, corpora-based methods in language engineering, electronic language resources and tools, formalization of natural languages

Digital Object Identifier 10.4230/OASICS.LDK.2021.10

Supplementary Material

Dataset: <http://redac.univ-tlse2.fr/lexiques/treelexPlusPlus.html>

1 Introduction

For Natural Language Processing (e.g., Information Extraction, Syntactic Parsing, Text Generation), as well as language-oriented Digital Humanities applications (e.g., Discourse Analysis, stylometry), machine-tractable as well as human-readable large-scale lexical resources are still a very valuable asset, even in a scene which appears today dominated by robust Machine-Learning algorithms and giga-word corpora. For instance, even though syntactic parsing has seen great advances in the past 10 years, thanks to the development of Treebanks and dependency-annotated corpora, even the best parser fails to capture in a consistent and predictable way such an intuitive linguistic notion as transitivity. In this sense, (semi-)manually constructed lexicons are an indispensable complementary resource to corpus-driven resources (e.g., “word embeddings”, n-grams datasets). We see the symbolic/-Machine Learning divide as a consequence of the fact that each type of resource addresses a portion of the problem. Thus, the challenge contemporary NLP systems are facing today is more how to integrate different knowledge sources than to prove that one source is better – or more consistent – than the other. In this paper, we present TreeLex++, an extension of TreeLex [16], a syntactic lexicon for French, based on the French Treebank (FTB), enriched here with aspectual information. Different lexical resources have been devised over several decades for the automatic processing of French texts, in different theoretical frameworks: from the manually-encoded Lexicon-Grammar tables [13] framed in a distributionalist framework, to contemporary large-scale, semi-automatically induced lexicons such as the Leff [24, 23], or resources acquired by way of “serious games”, such as Jeux de Mots [17, 18]. Most of those



© Anna Kupść, Pauline Haas, Rafael Marín, and Antonio Balvet;
licensed under Creative Commons License CC-BY 4.0

3rd Conference on Language, Data and Knowledge (LDK 2021).

Editors: Dagmar Gromann, Gilles Sérasset, Thierry Declerck, John P. McCrae, Jorge Gracia, Julia Bosque-Gil, Fernando Bobillo, and Barbara Heinisch; Article No. 10; pp. 10:1–10:12



Open Access Series in Informatics

Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

lexical resources have focused on providing a formalized description of the main syntactic categories, with an emphasis on verbal predicates. In extending TreeLex with aspectual information, our goal is primarily to set up a large-scale aspectual characterization process of verbs. Secondly, we wish to provide the NLP and DH communities with a resource which combines corpus-induced syntactic characterizations¹ as well as basic aspectual distinctions, based on Vendler’s classification [25].

In the first sections, we present how TreeLex++ derives from the original FTB-induced TreeLex resource (Section 2 and 3). Then we move on to the presentation of our aspectual semantics characterization process (Section 4). In Section 5 we give a general overview of the present state of the resource. Section 6 is dedicated to conclusions and perspectives.

2 TreeLex

TreeLex is a syntactic lexicon automatically extracted from the French Treebank [1]. The lexicon contains ca. 2000 contemporary French verbs with their syntactic realizations and frequencies found in the FTB. The FTB is a corpus of newspaper texts (*Le Monde* newspaper, 1990–1993), in which constituent trees were originally encoded in XML format. In addition to lexical information for every word (category, lemma, person, number, gender etc.), the corpus provides a syntactic structure for each sentence: both syntactic groups and functions are indicated (see Figure 1).

```

<SENT argument="ETR" author="MINANGOY ROBERT" date="1990-01-19" nb="1006" textID="456">
  <NP fct="SUJ">
    <w cat="D" ee="D-def-ms" ei="Dms" lemma="le" mph="ms" subcat="def">Le</w>
    <w cat="A" ee="A-ord-ms" ei="Ams" lemma="deuxième" mph="ms" subcat="ord">deuxième</w>
    <w cat="N" ee="N-C-ms" ei="NCms" lemma="problème" mph="ms" subcat="C">problème</w>
  </NP>
  <VN>
    <w cat="V" ee="V--P3s" ei="VP3s" lemma="être" mph="P3s" subcat="">est</w>
  </VN>
  <NP fct="ATS">
    <w cat="D" ee="D-def-fs" ei="Dfs" lemma="le" mph="fs" subcat="def">la</w>
    <w cat="N" ee="N-C-fs" ei="NCfs" lemma="nourriture" mph="fs" subcat="C">nourriture</w>
  </NP>
  <w cat="PONCT" ee="PONCT-S" ei="PONCTS" lemma="." subcat="S">.</w>
</SENT>

```

■ **Figure 1** A sample of FTB sentence annotation.

The XML-based annotation schema has since been complemented with a more straightforward tabulated format, following the CoNLL specifications that were widely adopted after the CoNLL shared task on dependency-parsing [20].

The FTB annotation schema is centered around the verbal nucleus (VN) which makes syntactic dependents easily accessible. This corpus organization is exploited by [16] in order to obtain obligatory arguments and provide syntactic frames for verbs present in the FTB. The resulting lexicon, called TreeLex², provides a rich syntactic representation of each argument since both functions and their phrasal realizations are encoded. Example 1 shows a lexical entry for the transitive verb *entraver* ‘to impede’ which takes a nominal subject (SUJ:NP) and a nominal direct object (OBJ:NP).

¹ As opposed to theory-driven ones.

² http://redac.univ-tlse2.fr/lexiques/treelex_en.html

■ **Table 1** TreeLex functions with syntactic realizations.

Tag	Function	Possible phrasal realizations
SUJ	subject	NP, VPinf, Ssub
OBJ	direct object	NP, VPinf, Ssub
A-OBJ	indirect object introduced by <i>à</i>	VPinf, PP
DE-OBJ	indirect object introduced by <i>de</i>	VPinf, PP
P-OBJ	indirect prepositional object (other than <i>de</i> and <i>à</i>)	PP
ATS	subject complement	AP, NP, VPpart, VPinf, Ssub
ATO	direct object complement	AP, NP, VPpart, VPinf, Ssub
ref	obligatory reflexive clitic pronoun	CL
obj	other obligatory clitic pronoun	en, y

1. *entraver*: SUJ:NP,OBJ:NP

In Treelex, names of functions and syntactic constituents are adopted directly from the FTB notation, with two additions (**ref** and **obj**) for obligatory clitics, cf. Table 1. Arguments with clitic realizations are used to indicate reflexive verbs (ex., *se réjouir* ‘to rejoice’: SUJ:NP,ref:CL), idiomatic expressions (ex., *s’en sortir* ‘to cope/get through’: SUJ:NP,obj:en,ref:CL) or an impersonal subject (ex., *falloir* ‘to have to’: SUJ:il,OBJ:VPinf).

If a verb allows for different syntactic combinations (i.e., either a list of functions or different realizations), every frame is listed separately. Therefore, a single verb (more precisely, its lemma) can be found several times in the lexicon, see (2). As no semantic disambiguation was performed, this strategy aims at distinguishing potentially different senses associated with each frame. Here, in (2a-b), *voler* has the meaning of ‘to steal’ whereas in (2c) it can be translated as ‘to fly’.

2. (a) *voler*: SUJ:NP,OBJ:NP,A-OBJ:NP

(b) *voler*: SUJ:NP,DE-OBJ:NP

(c) *voler*: SUJ:NP

As noted on TreeLex’s website, an optional realization of specific arguments has been added manually, cf. (3).

3. *détruire*: SUJ:NP,(OBJ:NP)

Finally, since multi-word units are indicated in the FTB, TreeLex lists 465 multi-word verbs, such as *courir le risque* ‘to take a risk’ or *donner lieu* ‘to result/take place’.

3 Beyond TreeLex: towards TreeLex++

TreeLex contains 1912 verbs and 3229 entries, i.e., verb-frame couples, which correspond to 24660 verb occurrences³ attested in the FTB corpus. The resource provides a rich set of syntactic information and, as stated in [16, p.38], it can be easily integrated with other resources for NLP tasks such as parsing, or text generation. However, its relatively small size makes open-domain applications problematic.

³ We present here figures from the on-line TreeLex version, http://redac.univ-tlse2.fr/lexiques/treelex/treelex_verbs.csv.

On the other hand, TreeLex's size makes an in-depth qualitative linguistic study feasible. For example, it could be extended with semantic information to investigate interactions between semantic and syntactic properties of verbs. For French, several projects have produced lexical resources containing syntactic and semantic verbal properties, or different levels of semantic information, e.g., verbal semantic classes (LVF, cf. [10]), thematic roles (French FrameNet, cf. [7]) or lexical aspect (Nomage, cf. [3] or [9]). In the current project, we decided to focus on high-level syntax-semantics relationships and thus we augmented the syntactic frames in TreeLex with manually encoded aspectual information. Our approach differs from [3] or [9], as verbal aspect assignment is guided by corpus examples rather than by elicited sentences.⁴ Similarly to [9], aspect is assigned to a verb-frame couple rather than to a verb alone. Nevertheless, the level of detail of our aspectual classes is distinct both from [3] and [9]: we use only the four major Vendlerian classes⁵.

In order to prepare the TreeLex data for aspect assignment, several modifications have been adopted. First, all frames had to be represented in a uniform way. Therefore all syntactic arguments, whether optional or not, have been treated equally and indications of optional realizations have been removed. In particular, verbs such as *détruire* 'to destroy' in (3) were transformed into (4):

4. *détruire*: SUJ:NP,OBJ:NP

Second, we had to address the ambiguity in TreeLex entries. As shown in (2), TreeLex verbs may appear with several frames. According to [16], this affects about 40% of TreeLex verbs. Such multiple frames may indicate a polysemous and/or a polyaspectual verb. However, all different syntactic realizations of a single argument structure (the same sequence of functions) are listed as separate frames in TreeLex, see (5). This representation is therefore unclear: it may show a true semantic (meaning) difference or introduce an artificial syntactic (frame) ambiguity. For example, the direct object (OBJ) of the verb *déplorer* 'to regret/deplore' in (5) has two syntactic realizations (a nominal phrase, NP, or a subordinate phrase, Ssub) but this syntactic variation does not imply a difference in meaning.

5. (a) *déplorer*: SUJ:NP,OBJ:Ssub

(b) *déplorer*: SUJ:NP,OBJ:NP

In order to avoid such an artificial ambiguity, we grouped all frames which differed only by their phrasal realization. Therefore, the double nature of OBJ in (5) is currently represented as in (6).

6. *déplorer*: SUJ:NP,OBJ:NP/Ssub

In an effort to reduce semantic ambiguity, we decided to consider only verbs which, after syntactic grouping, appeared with a single syntactic frame. As a consequence, verbs such as *voler* in (2) have been excluded.⁶ Multi-word verbal units have been omitted as well, as their meaning is usually idiosyncratic and conventional. Moreover, due to their idiomatic nature, syntactic construction appears heavily constrained.

Finally, all remaining 1161 verbs have been coupled with examples extracted from the FTB. We collected corpus examples in order to illustrate how each frame is instantiated and to provide a real context for aspect assignment.

⁴ [3] use corpus examples to assign aspectual properties only to nouns. Verbs are annotated with no explicit contextual information.

⁵ See Section 4 for details.

⁶ This strategy does not replace a real semantic disambiguation since verbs which allow for a single syntactic frame may still be polysemous. This issue will be addressed in further sections.

■ **Table 2** The four situation types, based on [25].

Class	Dynamic	Durative	Telic
STATE	–	+	–
ACT	+	+	–
ACC	+	+	+
ACH	+	–	+

4 Incorporating lexical aspect

Aspectual information has been added manually to TreeLex verbs. Unlike grammatical aspect, lexical aspect refers to inherent semantic properties indicating the way in which predicates are structured in relation to time. In the most general terms, the properties in question have to do with the presence (or lack thereof) of an end point (limit or boundary), duration or dynamicity in the lexical structure of certain classes of verbs. Thus, for instance, the presence of a limit distinguishes between **telic** (i.e., a time-limited situation) and **atelic** verbs.

These semantic properties give rise to four major aspectual classes (cf. [25]): STATE, ACTIVITY (ACT), ACCOMPLISHMENT (ACC) and ACHIEVEMENT (ACH). Their semantic features are listed in Table 2.

4.1 Annotation procedure

Aspectual assignment is a relatively new task, in the field of natural language annotation. The research exposed here is therefore to be seen as the first steps towards a full-fledged syntactic/semantic lexical resource. Our aspect assignment procedure consisted in a double manual annotation by two experts in semantics. Our annotation procedure is therefore not a “standard” annotation process, since, after the initial annotation phase, a final adjudication phase took place in order to arrive at the annotations presented in the current version of Treelex++. This process, which departs from established annotation approaches, is to be considered as a way of ensuring consistency in the current phase, where aspectual tagging is entirely performed manually. Each verb has been considered along with its syntactic frame and the corresponding examples found in the FTB. The assignment task consisted in choosing one of the four classes (tags) in Table 2. Each decision was made after applying the usual tests presented in the literature on verb lexical aspect (see [12, 15, 25, 8, 27, 19, 6, 22], among others). We have used the following six tests (cf. Table 3):

- **T1**: progressive form of *être en train de* ‘to be V-ing’
- **T2**: question related to dynamicity *Que s’est-il passé hier?* ‘What happened yesterday?’
- **T3**: use of aspectual semi-auxiliaries *commencer à* ‘to start doing something’, *continuer de* ‘to keep on doing something’, *arrêter de* ‘to stop doing something’
- **T4**: duration complement *en x temps* ‘in x time’
- **T5**: duration complement *pendant x temps* ‘during x time’
- **T6**: imperfective paradox *V[temps inaccompli] IMPLIQUE V [temps accompli]* ‘V[imperfect tense] IMPLIES V [perfect tense]’

10:6 Enriching a Lexical Resource for French Verbs with Aspectual Information

■ **Table 3** A grid for the allocation of aspectual classes to TreeLex verbs.

Situation type	T1	T2	T3	T4	T5	T6
STATE	no	no	no	yes no	yes no	yes
ACT	yes	yes	yes	no	yes	yes
ACC	yes	yes	yes	yes	yes no	no
ACH	no	yes	no	no	no	no

In order to illustrate our procedure, let us take the verb *invoquer* ‘to invoke’ in one of the sentences where it appears in the corpus:

7. Pour justifier cette décision, la direction invoque la déprime du marché automobile.
 ‘To justify this decision, the management invokes the depression of the automobile market.’

- **T1:** This verb cannot appear in a progressive form: **La direction est en train d’invoquer la déprime du marché automobile.*
- **T2:** *La direction a invoqué la déprime du marché automobile* is an acceptable answer to the question *Que s’est-il passé hier?*
- **T3:** This verb cannot appear as a complement of *commencer*, *continuer*, etc.: **La direction a commencé/continué à invoquer la déprime du marché automobile.*
- **T4:** *invoquer* is not compatible with *en x temps*: **La direction a invoqué la déprime du marché automobile en deux heures.*
- **T5:** the sentence is not compatible with *pendant x temps* either: **La direction a invoqué la déprime du marché automobile en deux heures.* This sentence is only acceptable in an iterative reading.
- **T6:** *La direction invoquait la déprime du marché automobile* does not imply *La direction a invoqué la déprime du marché automobile.*

Thus, according to the battery of tests summarized in Table 4, *invoquer* in (7) should be assigned to the ACHIEVEMENT class.

■ **Table 4** Test results for (7).

	T1	T2	T3	T4	T5	T6
<i>invoquer</i>	no	yes	no	no	no	no

It is important to mention that verbs were annotated according to their meaning in the sentences found in the FTB corpus. Verbal polysemy was addressed only if different meanings appeared in the corpus. It is known that phrasal context can influence the verbal aspect ([8, 26] *inter alia*). Upon applying the tests presented above, plural subjects and direct objects were transformed into their singular forms, so as to avoid the effect that plural arguments can turn ACC predicates (*écrire un article en dix jours* ‘to write a paper in ten days’) into ACT ones (*écrire des articles pendant dix jours* ‘to write papers for ten days’). Likewise, we have used past perfective tenses (*Elle a travaillé (hier)* ‘She has worked (yesterday)’) in order to avoid a habitual reading which is usually obtained in imperfective senses (*Elle (travaillait/travaille) à la poste* ‘She (worked/works) at the post office’). Since imperfective tenses favour a habitual reading, the dynamicity property [\pm dynamic] of the verb becomes inaccessible. For similar reasons, frequency adverbs triggering iterative or habitual readings (*souvent* ‘often’, *tous les jours* ‘every day’) were not taken into account either, since they interfere with verbal aspectual features.

We obtain an aspectual characterization limited to the meanings appearing in the corpus. It is not an annotation of the verbs as lemmas, neither verbs in sentences, but rather an annotation of verbal structures (verb + arguments) in a discursive context, which allowed us to identify verbal meaning and to avoid polysemy as much as possible.

4.2 Annotation consistency assessment: Inter-Rater Reliability

Based on the annotation process outlined above, we have been able to estimate the inter-rater reliability (IRR), by taking into account the annotations produced by two annotators on 1161 verbs. The annotators are both experts in aspectual semantics. Each verb in the list has been annotated independently by each annotator, even though a final adjudication step yielded the annotations visible in the current version of the lexicon. Comparing the annotations produced by both annotators was necessary, in order to arrive at a consistent decision in the final resource. For example, *atteler* ‘to tie’ was initially labelled “ACT” by annotator 1, while annotator 2 was not sure of his annotation. After the first annotation phase, both annotators agreed to tag the entry as “ACT”. Therefore, for the purpose of assessing the inter-rater agreement, we consider the initial annotation, which counts as a disagreement case. Conversely, for *cerner* ‘to surround’, annotator 1 was not sure of her annotation, while annotator 2 initially labelled the entry as “ACH”. After confronting their annotations, both annotators finally agreed on labelling this entry as “ACC”. Again, this case counts as a disagreement between both annotators. As can be seen, the final decision does not reflect either annotator’s initial decision, which underlines the fact that aspectual annotation is a complex task. Cases such as the one discussed here therefore strongly advocate in favor of a post-annotation adjudication phase.

The following IRR statistics were produced using R packages: `{irr}`⁷ and `{irrCAC}`⁸. Assessing IRR is not a straightforward task, since many methods have been presented in the literature⁹. We choose to present “standard” IRR statistics, such as Cohen’s Kappa [5], in this preliminary stage, alongside Gwet’s “Agreement Coefficient” score AC1 [14, 28]. Since the present lexical resource is still under construction, these IRR scores are essentially a way of assessing the complexity of the aspectual annotation task presented here, and therefore the consistency of the annotation procedure. In the annotation task under consideration, each annotator had to categorize 1161 verbal entries into 4 major classes: ACC, ACH, ACT, STATE. In total, 3 hybrid classes were also considered, such as: ACC/ACH, ACH/ACT and STATE/ACH. For example *varier* ‘to vary’ was initially labelled “ACH/ACT” by annotator 1 (final decision: “ACT”). Finally, a “not sure” tag was also used. As a consequence, the initial list of verbal entries has been associated with 8 different tags, including “not sure”.

As can be seen in Table 5, both annotators agree on 82.6% of the cases, with an estimated 9.7% of chance agreement. The reported Kappa score (0.744) indicates a moderate inter-rater agreement¹⁰, which is not uncommon for complex tasks. In our case, this score can be largely attributed to the fact that 4 major classes and 3 hybrid ones were considered. Gwet’s AC1 score (0.806) is slightly higher than Cohen’s Kappa, which can be attributed to

⁷ Version 0.84.1, see [11] for more details on the underlying implementation, and [21] for a presentation of the R platform.

⁸ Version 1.0, see [14] for a comprehensive presentation of the chance-corrected agreement coefficients implemented in this package.

⁹ See [2] for a survey of IRR methods in NLP.

¹⁰ Assessing the relevance of Kappa scores is known to depend heavily on the domain of application. We see these scores as an estimation of the task’s complexity as well as the overall quality of the proposed annotations.

■ **Table 5** IRR assessment of TreLex++ aspectual annotations.

Method	Score
irr (2 raters)	
unweighted Cohen’s Kappa	0.744
irrCAC (confidence level = 0.95)	
percent agreement pa	0.826
percent chance agreement pe	0.097
AC1	0.806

the fact that Gwet’s AC1 is a chance-corrected agreement coefficient that is known to yield higher agreement coefficients than Cohen’s (and other authors’), in certain configurations. Regardless of the method, these figures indicate a “moderate” to “good” inter-annotator agreement.

At this point, it is worth emphasizing once more that, once the preliminary annotation was completed, a final adjudication phase took place, which yielded the final aspectual annotations visible in the current version of Treelex++. Since these final annotations are those end users will see, it is necessary to assess IRR scores between each annotator and the final annotations. In this case, Kappa scores in the 0.85 range, and AC1 scores in the 0.9 range can be reported. Final users of the TreeLex++ lexical resource should therefore consider that the proposed aspectual annotations are consistent, and that the annotation procedure based on syntactico-semantic tests achieves good results for the classes considered. As encouraging as they might seem, these figures should not obscure the fact that there is still considerable room for improvement, in terms of both scale and detail. For future versions, we are contemplating Games With A Purpose (GWAP) such as JeuxdeMots [17] as a source of user input. We are confident JeuxdeMots players will consider favorably new games, such as aspect-oriented tasks, provided we are able to propose ‘gamified’ versions of the present annotation procedure.

5 Data in TreeLex++

The resulting resource, TreeLex++, contains 1161 verbs enriched with syntactic (frame) and semantic (lexical aspect) properties. It is available in a text format as a CSV file (comma separated value). Each verb is accompanied by its frame, the lexical aspect, the number of examples found in the FTB and their full list¹¹. To simplify the search of the inflected form in the example text, the corresponding verb is indicated between and tags, as presented in (8):

8. Quant à moi , je trouve qu’ on se fiche du monde en n’ expliquant pas les choses en langage courant .
 ‘As for me, I think that they don’t give a toss about the people by giving no explanation in the common language.’

To make linguistic generalizations easier, information encoded in syntactic frames has been translated into several representations:

¹¹ Individual examples are separated by a vertical bar ‘|’.

- number of syntactic arguments¹²
- whether a verb is reflexive or not
- a general frame (a list of syntactic functions and obligatory clitics)
- a simplified frame (a list of syntactic functions alone)
- the full frame including syntactic realizations (types of phrases)

The corresponding syntactic information for *déplorer* in (6) and the reflexive verb *se ficher* ‘to not give a toss’ presented in TreeLex++ format is given in Table 6.

■ **Table 6** Syntactic information in TreeLex++.

Verb	Number of Arguments	Reflexive?	General frame	Simplified frame	Full frame
<i>déplorer</i>	2 2	no	SUJ.OBJ	SUJ.OBJ	SUJ:NP.OBJ: NP/Ssub
<i>se ficher</i>	2	yes	SUJ.DE-OBJ.refl	SUJ.DE-OBJ	SUJ:NP.DE-OBJ .refl:CL

A brief summary of syntactic realizations¹³ of TreeLex++ verbs is given in Table 7 below. The number of arguments in TreeLex++ does not exceed three and the vast majority of verbs (74.24%) have two arguments. However, as indicated in Table 6, this does not necessarily correspond to a transitive structure (SUJ.OBJ) as the second argument may have a different function than a direct object (see Table 1).

■ **Table 7** The distribution of verbs with respect to the number of arguments.

Number of Arguments	Total	Percent
1	183	15.76%
2	862	74.24 %
3	116	9.99%

The distribution of verbal aspectual classes found in TreeLex++ is given in Table 8.

■ **Table 8** Aspect distribution in TreeLex++.

Aspectual class	Total	Percent
ACH	576	49.61%
ACC	260	22.39%
ACT	219	18.86%
STATE	103	8.87%
polysemous verbs	3	0.27%

The majority of verbs in TreeLex++ are telic (ACH or ACC). If we look at dynamicity, only a small proportion of verbs (8.87%) are true statives, the bulk of the entries are dynamic (ACH, ACC or ACT). However, the distribution of durative (STATE, ACT, ACC) and non-durative (ACH) verbs is almost equal.

¹² Clitic arguments are not considered here.

¹³ The number of syntactic arguments.

10:10 Enriching a Lexical Resource for French Verbs with Aspectual Information

The resource is neither syntactically nor semantically balanced, which is probably due to the content of the FTB corpus (newspaper texts).

As shown in Table 8, most verbs are assigned a single aspect. Hence, it seems that our approximate disambiguation technique is quite efficient. 3 verbs, however, exhibit a double aspect: *excéder*, *observer*, and *traverser*. Indeed, judging from their context, these verbs are truly polysemous in the FTB: *excéder* is ambiguous between ‘to exceed’ and ‘to infuriate’, *observer* is used as either ‘to observe’ or ‘to respect/keep’ and *traverser* corresponds to ‘to cross’ or ‘to experience’. Therefore, even when syntactic properties are restricted to a single frame, certain semantic ambiguities could remain.

6 Conclusions and perspectives

TreeLex++ is a lexical resource which associates both syntactic and semantic properties, for over a thousand verbs, illustrated with attested examples taken from the FTB. Such a database offers a valuable resource for fundamental linguistics research, NLP and DH applications. From a fundamental research perspective, TreeLex++ allows to identify correlations, if any, between syntactic frames and aspect values. In other words, it allows researchers to work at the syntax/semantics interface. For instance, intuitively, the accomplishment verbs (ACC) should be associated with transitive verbs (2-argument predicates). TreeLex++ provides an opportunity to verify this hypothesis empirically: not only can it be confirmed or refuted but we can also estimate the degree of association between syntactic structures and aspectual classes. The first findings presented in [4] show how TreeLex++ can be put to use in this perspective. As for NLP applications, a number of practical uses of aspectual information is cited in [9]: the assessment of event factuality, text summarization, machine translation or automatic detection of temporal relations. We anticipate performance gains for those task, by integrating TreeLex++ as a symbolic resource, within a Machine Learning processing chain.

In its current version, TreeLex++ contains only single-frame verbs, which roughly covers a half of the entries in TreeLex. In order to include the remaining half in TreeLex++, we have to employ a true semantic disambiguation technique first. As mentioned in Section 5, a verb with a unique syntactic combination may still be polysemous and polyaspectual. In case of several frames, this potential ambiguity is multiplied and human disambiguation effort, already complex and time-consuming, increases considerably. A possible solution could be a lexical look-up of verb-frame couples in LVF [10] in order to identify different verb senses. However, pairing the senses with the corresponding FTB examples would require an ad-hoc approach. As mentioned above, another available option is to leverage user input, by resorting to crowd-sourcing, or “Game With A Purpose” platforms. We have taken steps towards this end by contacting Jeux de Mots’s developer, Mathieu Lafourcade, in the perspective of integrating the aspectual information from TreeLex++ to the existing Jeux de Mots lexical network. This will allow for the development of new types of lexical games. We also hope Lafourcade’s lexical propagation and integrity checking mechanisms will allow us to capture more general syntax/semantics properties than those which can be currently found in the FTB.

An evaluation methodology for our resource is also in order, beyond Inter-Rater Reliability scores, to determine the accuracy, as well as the coverage of our aspectual assignment process. For instance, we could compare our results with aspect values attributed to verbs in the Nomage project [3]. However, Nomage methodology (for verbs) differs from ours as aspect assignment is based on elicited examples rather than on verb uses in a corpus. Another

comparison could be made with the syntactico-semantic resource described in [9] which served for training of an automatic classifier of verbal aspect. Unfortunately, this data does not seem to be publicly available. Moreover, both resources use different aspectual values from ours thus the corresponding tagsets have to be converted first in order to provide the equivalent information. Again, we turn towards the Jeux de Mots platform, in the hope of gaining insights from users’s inputs on lexical aspect assignment tasks¹⁴, as well as from the network’s built-in sanity checking mechanisms.

The current version of TreeLex++ is freely available on-line: <http://redac.univ-tlse2.fr/lexiques/treelexPlusPlus.html>. It can be either downloaded as a text (CSV) file or browsed directly via an intuitive on-line interface: <http://redac.univ-tlse2.fr/lexiques/treelexPlusPlus/interface/TreelexPlusPlusBrowser.html>.

References

- 1 Anne Abeillé, Lionel Clément, and François Toussnel. Building a treebank for French. In *Treebanks*, pages 165–187. Springer, 2003.
- 2 Ron Artstein and Massimo Poesio. Inter-coder agreement for computational linguistics. *Computational Linguistics*, 34(4):555–596, 2008.
- 3 Antonio Balvet, Lucie Barque, Marie Hélène Condette, Pauline Haas, Richard Huyghe, Rafael Marin, and Aurélie Merlo. La ressource Nomage. Confronter les attentes théoriques aux observations du comportement linguistique des nominalisations en corpus. *Traitement Automatique des Langues*, 52(3):129–152, 2011.
- 4 Antonio Balvet, Pauline Haas, Anna Kupść, and Rafael Marin. Looking for Syntax/Aspect Mappings: a Case Study on the French Treebank. In *Grammar and Corpora*, 2018.
- 5 Jacob Cohen. A coefficient of agreement for nominal scales. *Educational and psychological measurement*, 20(1):37–46, 1960.
- 6 Anne Daladier. Le rôle des verbes supports dans un système de conjugaison nominale et l’existence d’une voix nominale en français. *Langages*, pages 35–53, 1996.
- 7 Marianne Djemaa, Marie Candito, Philippe Muller, and Laure Vieu. Corpus annotation within the French Framenet: a domain-by-domain methodology. In *Tenth International Conference on Language Resources and Evaluation (LREC 2016)*, 2016.
- 8 David R Dowty. *Word Meaning and Montague Grammar: The Semantics of Verbs and Times in Generative Semantics and in Montague’s Ptq*, volume 7. Springer Science & Business Media, 1979.
- 9 Ingrid Falk and Fabienne Martin. Automatic identification of aspectual classes across verbal readings. In *Sem 2016 The Fifth Joint Conference on Lexical and Computational Semantics*, 2016.
- 10 Jacques François, Denis Le Pesant, and Danielle Leeman. Présentation de la classification des verbes français de Jean Dubois et Françoise Dubois-Charlier. *Langue française*, 1(153):3–19, 2007. doi:10.3917/lf.153.0003.
- 11 Matthias Gamer, Jim Lemon, and Ian Fellows Puspendra Singh. *irr: Various Coefficients of Interrater Reliability and Agreement*, 2019. R package version 0.84.1. URL: <https://CRAN.R-project.org/package=irr>.
- 12 Howard B Garey. Verbal aspect in French. *Language*, 33(2):91–110, 1957.
- 13 Maurice Gross. *Méthodes en syntaxe*. Hermann, Paris, 1975.
- 14 Kilem L Gwet. *Handbook of inter-rater reliability: The definitive guide to measuring the extent of agreement among raters*. Advanced Analytics, LLC, 2014.
- 15 Anthony Kenny. *Action, emotion and will*. Routledge, 2003.

¹⁴In our view, the set of tests presented in Section 3 could very well be adapted to new games, focusing on aspectual properties.

10:12 Enriching a Lexical Resource for French Verbs with Aspectual Information

- 16 Anna Kupść and Anne Abeillé. Growing Treelex. In *International Conference on Intelligent Text Processing and Computational Linguistics*, pages 28–39. Springer, 2008.
- 17 Mathieu Lafourcade. Making people play for Lexical Acquisition with the JeuxDeMots prototype. In *SNLP'07: 7th International Symposium on Natural Language Processing*, page 7, Pattaya, Chonburi, Thailand, 2007. URL: <https://hal-lirmm.ccsd.cnrs.fr/lirmm-00200883>.
- 18 Mathieu Lafourcade and Alain Joubert. Détermination des sens d'usage dans un réseau lexical construit grâce à un jeu en ligne. In *TALN'08: Traitement Automatique des Langues Naturelles*, pages 189–199, 2008.
- 19 Béatrice Lamiroy. The complementation of aspectual verbs in French. *Language*, pages 278–298, 1987.
- 20 Joakim Nivre, Johan Hall, Sandra Kübler, Ryan McDonald, Jens Nilsson, Sebastian Riedel, and Deniz Yuret. The CoNLL 2007 shared task on dependency parsing. In *Proceedings of the 2007 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning (EMNLP-CoNLL)*, pages 915–932, 2007.
- 21 R Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria, 2013. URL: <http://www.R-project.org/>.
- 22 Susan Rothstein. *Structuring events: A study in the semantics of lexical aspect*, volume 5. John Wiley & Sons, 2008.
- 23 Benoît Sagot. The Lefff, a freely available and large-coverage morphological and syntactic lexicon for French. In *7th international conference on Language Resources and Evaluation (LREC 2010)*, 2010.
- 24 Benoît Sagot, Lionel Clément, Eric Villemonte de La Clergerie, and Pierre Boullier. The Lefff 2 syntactic lexicon for French: architecture, acquisition, use. In *LREC 06*, pages 1–4, 2006.
- 25 Zeno Vendler. *Linguistics in philosophy*. Cornell University Press, 1967.
- 26 Henk Verkuyl. A theory of aspectuality (cambridge studies in linguistics 64), 1993.
- 27 Marc Wilmet. Aspect grammatical, aspect sémantique, aspect lexical: un problème de limites. *J. David; R. Martin, (éds), La notion d'aspect*, Metz: Centre d'Analyse Syntaxique de l'Université de Metz, pages 51–68, 1980.
- 28 Nahathai Wongpakaran, Tinakon Wongpakaran, Danny Wedding, and Kilem L Gwet. A comparison of Cohen's Kappa and Gwet's AC1 when calculating inter-rater reliability coefficients: a study conducted with personality disorder samples. *BMC medical research methodology*, 13(1):1–7, 2013.