



Automated Semantic Classification of French Verbs

Ingrid Falk

► **To cite this version:**

Ingrid Falk. Automated Semantic Classification of French Verbs. Document and Text Processing. 2008. <hal-01075493>

HAL Id: hal-01075493

<https://hal.inria.fr/hal-01075493>

Submitted on 17 Oct 2014

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.

CRÉATION AUTOMATIQUE DE CLASSES SÉMANTIQUES VERBALES POUR LE FRANÇAIS

AUTOMATED SEMANTIC CLASSIFICATION
OF FRENCH VERBS

Rapport de Master Informatique Master Thesis in Computer Science

spécialité **Traitement Automatique des Langues**
Computational Linguistics

Auteur:
Ingrid FALK

Encadrants:
Claire GARDENT
Fabienne VENANT

Jury:
Noëlle CARBONNEL
Didier GALMICHE
Claude GODART
Dominique MÉRY
Guy PERRIER

— 27 juin 2008 —

Abstract

The aim of this work is to explore (semi-)automatic means to create a Levin-type classification of French verbs, suitable for Natural Language Processing. For English, a classification based on Levin's method ([Lev93]) is *VerbNet* ([Sch05]). *VerbNet* is an extensive digital verb lexicon which systematically extends Levin's classes while ensuring that class members have a common semantics and share a common set of syntactic frames and thematic roles.

In this work we reorganise the verbs from three French syntax lexicons, namely *Volem*, the Grammar-Lexicon (Ladl) and *Dicovalence*, into VerbNet-like verb classes using the technique of *Formal Concept Analysis*.

We automatically acquire syntactic-semantic verb class and diathesis alternation information. We create large scale verb classes and compare their verb and frame distributions to those of VerbNet.

We discuss possible evaluation schemes and finally focus on an evaluation methodology with respect to VerbNet, of which we present the theoretical motivation and analyse the feasibility on a small hand-built example.

Résumé

L'objectif de ce travail est d'explorer dans quelle mesure l'information contenue dans trois lexiques syntaxiques pour le Français (*Volem*, le *Lexique-Grammaire* et *Dicovalence*) peut être utilisée pour regrouper les verbes en classes sémantiques à la Beth Levin [Lev93]. Ce type de classification a été réalisé pour l'anglais avec *VerbNet* ([Sch05]). *VerbNet* est un lexique verbal électronique qui reprend la classification de Beth Levin et l'étend systématiquement en assurant la cohérence sémantique et syntaxique de ses classes.

Dans ce travail nous utilisons la technique de *Analyse Formelle de Concepts* pour réorganiser les verbes des trois lexiques syntaxiques pour le Français en classes verbales à la Beth Levin.

Nous montrons comment l'Analyse Formelle de Concepts permet la création automatique de classes verbales à grande échelle et nous comparons la répartition de verbes et cadres syntaxiques dans ces classes à celle de VerbNet.

Deuxièmement, l'Analyse Formelle de Concepts nous permet de découvrir les paires de constructions verbales les plus fréquentes présentes dans les lexiques. Ces paires de cadres syntaxiques représentent des candidats d'alternances verbales, qui sont fondamentales à la réalisation d'une classification à la Beth Levin.

Nous discutons plusieurs scénarios d'évaluations pour nous concentrer finalement sur une méthode d'évaluation par rapport à VerbNet, dont nous présentons les motivations théoriques et analysons la faisabilité sur un petit exemple.

Contents

| | | |
|----------|---|-----------|
| 1 | Introduction | 1 |
| 2 | Verb classes | 2 |
| 2.1 | Verb classes for English | 2 |
| 2.1.1 | Beth Levin's Verb Classes | 2 |
| 2.1.2 | VerbNet | 3 |
| 2.1.3 | Why are Verb Classes Useful? | 5 |
| 2.2 | Verb Classes for French | 5 |
| 3 | Goals and Methods | 6 |
| 4 | The Resources | 8 |
| 4.1 | Volem | 8 |
| 4.2 | The Ladi Grammar-Lexicon | 9 |
| 4.3 | Dicovalence | 9 |
| 4.4 | Coverage Summary | 10 |
| 5 | Formal Concept Analysis | 10 |
| 5.1 | Formal concepts | 11 |
| 5.2 | Association Rules | 13 |
| 6 | Extraction of Verb Class Information from Syntactic Lexicons | 14 |
| 6.1 | Building Verb Classes | 15 |
| 6.2 | Mining Diathesis Alternations for French | 20 |
| 6.3 | Discussion | 21 |
| 6.3.1 | Evaluation | 22 |
| 6.3.2 | Conclusion and Direction for further Work | 23 |
| 7 | Possible Evaluation Methodology | 24 |
| 7.1 | Relating concepts using Ontology Based Concept Similarity | 24 |
| 7.2 | The Running Example | 27 |
| 7.3 | Concept Similarity for the Running Example | 29 |
| 7.4 | Comparing Class Hierarchies | 31 |
| 8 | Conclusion | 31 |
| A | Tables | 33 |
| B | Implementation | 36 |

1 Introduction

Verb classes categorise verbs into classes such that verbs in the same class are as similar as possible, and verbs in different classes are as dissimilar as possible. Several kinds of classifications exist: For example, syntactic verb classes categorise verbs according to syntactic properties, semantic verb classes categorise verbs according to semantic properties.

From a practical point of view, verb classes reduce redundancy in a verb lexicon, in that they list the common properties of the verbs belonging to the same classes. In addition, verb classes can help predict properties of a verb that received insufficient empirical evidence, by referring to verbs in the same class: thereby a verb classification is especially useful for the recurrent problem of data sparseness in Natural Language Processing (NLP), where little or no knowledge is available for rare events.

This work is concerned with the *automatic* or *semi-automatic* creation of verb classes for French at the *syntax-semantics interface*. Recently much work has concentrated on the acquisition of this type of classifications because semantic information from corpora is often not available. Instead, researchers typically make use of syntax-semantics verb classes by exploiting a long standing-linguistic hypothesis, namely that there is a tight connection between the lexical meaning of a verb and its behaviour. Even though the meaning-syntax relationship is not perfect, researchers use the assumption that a verb classification based on the syntactic behaviour of verbs agrees to a certain extent with a semantic classification.

The aim of this work is to explore ways towards an automatic or semi-automatic syntactic-semantic classification of French verbs.

For English, there are two important verb classifications of this type: the first is Beth Levin's work ([Lev93]), the second is VerbNet ([Sch05]), the largest *computational* verb lexicon currently available. The latter provides detailed syntactic-semantic descriptions of Levin's classes and organises them systematically into a refined hierarchy. Because of its coherence, coverage and its suitability for NLP it comes very close to the type of resource we would like to achieve for French.

For French, there exist several resources which contain rich and extensive information about the morphosyntactic and semantic properties of French verbs. They are organised in different ways and based on different assumptions, but they all provide information on single verbs rather than verb classes. Hence the objective of this work is to obtain verb *class* information from some of the syntactic-semantic properties assigned to single verbs by these resources. We worked with three of the French resources, namely Volem ([FVSD⁺02]), the Grammar-Lexicon ([Gro75], [GL92], [BGL76]) and Dicovalence (<http://bach.arts.kuleuven.be/dicovalence/>).

Formal Concept Analysis is one of many applicable classification and clustering techniques. This methodology seems well adapted to this task for several reasons: First, it operates on two types of items (called objects and attributes) related by a binary relation – in our approach verbs will be the objects and syntactic properties the attributes. Second, it yields classes defined by both a subset of objects, and a subset of attributes. Third, it organises the classes in a hierarchical structure which naturally comes close to that of VerbNet, that is the subsumption relations between the sets of class verbs on one hand and the sets of class properties on the other are similar.

After creating verb classes we need to judge their similarity to the targeted classification. To our knowledge there is no absolute or generally accepted evaluation method. However, we discuss the issues relating to evaluation raised by our results and go some way towards their systematic analysis.

In Section 2 we present the background for our work namely the research done on verb classes for English and for French. In Section 3 we give an overview of the outcome of this work and the approaches used to achieve it. Section 4 describes the French resources we use in more detail and Section 5 briefly introduces *Formal Concept Analysis* (FCA). In Section 6 we describe how FCA was applied to obtain verb class information and show how the results motivated the direction for further work, namely a particular evaluation scheme. Finally, in Section 7 we present the theoretical background for this evaluation methodology and test it on a small hand-built example.

2 Verb classes

2.1 Verb classes for English

Today's work on verb classes for English is strongly influenced by Beth Levin's seminal work ([Lev93]). In Section 2.1.1, we start by introducing her work. We then (Section 2.1.2) present VerbNet, the most extensive digital resource for English verb classes currently available. Finally, Section 2.1.3 summarises the motivations for the acquisition of verb classes.

2.1.1 Beth Levin's Verb Classes

In her "Preliminary Investigation", Beth Levin provides a classification of English verbs which is guided by the hypothesis that there is a systematic relation between the syntactic and semantic properties of verbs. More specifically

- (1) If the members of a set of verbs share some meaning component, then the members of this set can be expected to exhibit the same syntactic behaviour(s) and conversely,
- (2) If the members of a set of verbs exhibit the same syntactic behaviour(s), then the members of this set can be expected to share some meaning component(s).

In [Lev93], Levin systematically and consequently creates an important empirical basis verifying these ideas.

Her underlying research methodology is the following:

- (1) first she sets out to define a range of **diathesis alternations**.
- (2) then the diathesis alternations are used to isolate semantically coherent **classes of verbs**.

Diatheses alternations are variations in the (syntactic) realisation of a verb argument, sometimes accompanied by meaning changes.

For instance, the following are alternations taken from [Lev93] :

- (1) the MIDDLE alternation:

| | |
|------------------------|---------------|
| He cuts the bread. | $NP_0 V NP_1$ |
| The bread cuts easily. | $NP_1 V ADV$ |

- (2) the CONATIVE alternation:

| | |
|-----------------------|--------------------|
| He cuts the bread. | $NP_0 V NP_1$ |
| He cuts at the bread. | $NP_0 V PP_1 [at]$ |

(3) the BODY-PART POSSESSOR alternation:

He cut his arm. NP₀ V NP₁
 He cut himself on the arm. NP₀ V REFL₁ PP₂

The second column in the examples above shows how diathesis alternations can naturally be represented by *syntactic* argument frames¹.

Verb classes. Verb subgroups pattern together with respect to the set of diathesis alternations they can enter.

For example *cut* obviously is acceptable in all three alternations above, but *touch* only enters the BODY-PART POSSESSOR alternation:

* The bread touches easily MIDDLE
 * He touches at the bread CONATIVE
 She touches him on the arm BODY-PART POSSESSOR

The pattern of behaviour for the verbs *touch*, *hit*, *cut* and *break* is summarised in Table 1.

| | <i>touch</i> | <i>hit</i> | <i>cut</i> | <i>break</i> |
|---------------------|--------------|------------|------------|--------------|
| MIDDLE | No | No | Yes | Yes |
| CONATIVE | No | Yes | Yes | No |
| BODY-PART POSSESSOR | Yes | Yes | Yes | No |

Table 1: Pattern of behaviour with respect to three diathesis alternations

As other verbs pattern like each of these four verbs, these diathesis alternations induce four verb classes:

Break verbs: break, crack, rip, shatter, snap, ...
Cut verbs: cut, hack, saw, scratch, slash, ...
Touch verbs: pat, stroke, tickle, touch, ...
Hit verbs: bash, hit, kick, pound, tap, whack, ...

In this way, Beth Levin identified a set of 79 diathesis alternations and manually classified about 3200 English verbs in about 200 classes. Nonetheless, it is still incomplete and not available electronically, which limits its usefulness for NLP (Natural Language Processing).

2.1.2 VerbNet

VerbNet ([Sch05]) aims to produce an electronic verb lexicon that is based on Beth Levin's classification methodology and extensively covers English verbs. Currently, the largest on-line lexicon² for English verbs and verb classes, VerbNet extends Levin's approach as follows:

- It is a hierarchical and domain-independent extension of Levin's classification through refinement and addition of subclasses in order to achieve syntactic and semantic coherence,

¹A verb's argument frame defines the type of phrase constituents the verb may be combined with and the way these constituents pattern with the verb to build a grammatical construction. For instance, in the examples NP₀ V NP₁ designates an argument frame for the verb *cut*. It shows that *cut* can be used in a construction with a noun phrase (NP) as subject and direct object respectively.

²<http://verbs.colorado.edu/verb-index/index.php>

- broad-coverage: about 5200 verbs and 237 classes,
- It ensures that all members of each (sub)class have a common semantics and share a common set of syntactic frames and thematic roles.

VerbNet was created essentially manually taking on an initial classification based on a refinement of Levin’s classes. Due to the careful design of the lexicon it was then possible to extend it with various and rather heterogeneous other classifications, showing ways towards the use of semi-automatic methods. Moreover it was possible to establish mappings to other popular lexical resources such as WordNet³, Xtag⁴, and FrameNet⁵, making it into an important and novel knowledge base for NLP.

Each VerbNet class is completely described by

1. its set of members,
2. thematic roles,
3. syntactic frames and
4. selectional restrictions (optional).

A simplified example entry for the class *Hit-18.1* is given in Table 2.

Currently, the relevant items for our work are the verb members and the **syntactic frames**.

Syntactic frames describe constructions such as transitive, intransitive, prepositional phrases and most of the components of Levin’s alternations. They are represented by a brief description (eg. BASIC TRANSITIVE or CONATIVE in Table 2), an example, a syntactic description (eg. *Agent V at Patient*) and a set of semantic predicates (column **Semantics** in Table 2). The syntactic description consists of the thematic roles in their preferred argument position, the verb itself and other lexical items, eg. prepositions, required for particular constructions. In what follows, we will only make use of the frame **name** and **syntax**.

Each member of a class accepts all the syntactic frames of the class and must accept all the classes’ thematic roles in the corresponding argument slots. Verb classes are organised in an inheritance hierarchy. A class may be subdivided into subclasses, according to specific syntactic frames or thematic roles which are valid only for a subset of the class members. A verb with a specific sense can only belong to exactly one class. The information presented in the class is strictly monotonic: Each subclass adds more information in terms of syntactic frames or thematic roles or imposes further restrictions to its parent to the ones already present.

³WordNet® is a freely and publicly available (<http://wordnet.princeton.edu/>) large lexical database of English, where nouns, verbs, adjectives and adverbs are grouped into sets of cognitive synonyms (synsets). Synsets are interlinked by means of conceptual-semantic and lexical relations. Its structure and large coverage make it a useful tool for computational linguistics and NLP ([Fel98]).

⁴The Xtag project (<http://www.cis.upenn.edu/~xtag/>) developed a comprehensive English grammar based on the Tree Adjoining Grammar formalism together with a set of tools. They are public domain resources providing very detailed syntactic characterisations of English verbs which explains their primary usage as a large-scale parsing resource in NLP.

⁵The Berkeley FrameNet project (<http://framenet.icsi.berkeley.edu/>) is creating an on-line lexical resource for English, based on frame semantics and supported by corpus evidence. The aim is to document the range of semantic and syntactic combinatorial possibilities (valences) of each word in each of its senses, through computer-assisted annotation of example sentences.

| Class | | <i>Hit-18.1</i> | |
|---------------------------------|--|---------------------|---|
| Parent | – | | |
| Thematic roles | Agent, Patient, Instrument | | |
| Selectional restrictions | Agent[+int.control] Patient[+concrete] Instrument[+concrete] | | |
| Frames | | | |
| Name | Example | Syntax | Semantics |
| <i>Basic Transitive</i> | Paula hit the ball | Agent V Patient | cause(Agent, E) manner(during(E), directedmotion, Agent) not(contact(during(E), Agent, Patient)) manner(end(E), forceful, Agent) contact(end(E), Agent, Patient) |
| | Paula kicked the door open | Agent V Patient Adj | cause(Agent, E) manner(during(E), directedmotion, Agent) not(contact(during(E), Agent, Patient)) manner(end(E), forceful, Agent) contact(end(E), Agent, Patient) Pred(result(E), Agent, Patient) |
| <i>Resultative</i> | Paul hit at the window | Agent V at Patient | cause(Agent, E) manner(during(E), directedmotion, Agent) not(contact(during(E), Agent, Patient)) |
| <i>Conative</i> | | | |

Table 2: Simplified VerbNet entry for the *Hit-18.1* class

For instance, the class *Hit-18.1-1* is a subclass of *Hit-1.18*. It imposes the additional selectional restriction *body-part* or *reflexive* on the *Instrument* thematic role and adds the frame TRANSITIVE BODY-PART OR REFLEXIVE OBJECT as in *Paul hit her elbow*.

VerbNet clearly is a model for the type of resource we would like to achieve for French.

2.1.3 Why are Verb Classes Useful?

Beth Levin’s original motivation was mainly theoretical: she was interested in identifying primitives for lexical semantics and explored the idea that semantic differences were reflected in the syntax (verbs sharing alternations were verbs sharing some lexical semantic primitives).

From a more practical point of view however, verb classification brings a number of benefits. First, verb classes permit factorising a verb lexicon and thereby facilitate its extension, debugging and maintenance. Second, verb classes have turned out to be very helpful in a variety of Natural Language Processing tasks such as language generation, subcategorisation acquisition and document classification. Third, if verb classes share commonalities across languages (a hypothesis still to be verified), they might provide a powerful way to support multilingual systems.

2.2 Verb Classes for French

There exist several resources for French which contain rich and extensive information about the morphosyntactic and semantic properties of French verbs. In this thesis, we will focus on the following three:

Volem [FVSD⁺02], is the outcome of a regional European project where the aim was to design a lexical knowledge base for verbs, where syntactic and semantic descriptions are normalised and treated in a uniform way across several (Romance) languages. It is based on earlier work of Patrick Saint-Dizier's ([SD96]) where he attempts to build a French resource comparable to that of Beth Levin's.

The Grammar-Lexicon, aka **Ladl tables**⁶, [Gro75, BGL76, GL92] is a resource developed by Maurice Gross and his collaborators at the "Laboratoire d'Automatique Documentaire et Linguistique" from 1968 to 2002. It classifies verbs according to their syntactic properties and distributional constraints on the constructions they accept by assigning them to "tables".

Dicovalence, <http://bach.arts.kuleuven.be/dicovalence/>, is a syntactic lexicon for French assigning certain subcategorisation information, called "valency frames", to verbs. Initially developed at the University of Leuven (Belgium) from 1986 to 1992 within the Proton research project, it has recently been completely reviewed and updated under the direction of Karel van den Eynde and Piet Mertens.

Albeit constructed in different ways and based on different assumptions, each of these resources is in essence a syntactic lexicon which associates a verb with a description of the argument frames⁷ they allow. In other words, they differ from a verb classification such as VerbNet or Beth Levin's classes in that they fail to provide information about verb classes and/or verb alternations. In the following sections, we will present these resources in more detail (Section 4) and show how they can be used to help define a semi-automatic means of producing verb classes for French.

3 Goals and Methods

The overall goal – to obtain a VerbNet like classification for French exploiting the resources mentioned above – is way beyond the scope of a master thesis.

A more modest and accessible objective, which still goes some way in this direction is to explore the adequacy of well-known semi-automatic classification techniques for this purpose.

We selected the classification techniques based on *Formal Concept Analysis* because its principles – building clusters simultaneously on objects and on their properties – seemed particularly well suited to our problem (cf. Sec. 6.1).

Formal Concept Analysis is a method for deriving conceptual structures out of data. The process starts from a set of objects, a set of properties (or attributes) and a relation table (called *formal context*) defining the properties a given object has. FCA then establishes a partial order on concepts, where concepts are pairs of objects and attributes subsets.

This partial order can be graphically represented in a hierarchy, facilitating the analysis of complex structures and the discovery of dependencies within the data.

We will use FCA by considering the verbs as objects. Properties will be information items provided by the French resources introduced previously (cf. Section 2.2 and Section 4).

For example, consider the following information extracted from B. Levin's book:

⁶<http://ladl.univ-mlv.fr/DonneesLinguistiques/Lexiques-Grammaires/Presentation.html>

⁷A verb's argument frame defines the type of phrase constituents the verb may be combined with and the way these constituents pattern with the verb to build a grammatical construction. For instance *Agent V at Patient* in the VerbNet sample entry for the *Hit-18.1* class (Table 2) is an argument frame for *hit*. It shows that *hit* can be used in a construction like *Paul hit at the window*.

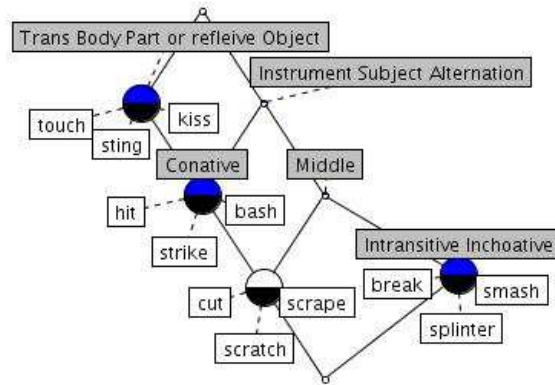


Figure 1: The concept lattice for the formal context 3.

- the verbs** break, splinter, smash, cut, scratch, scrape, hit, strike, bash, touch, sting, kiss and
- the alternations** MIDDLE, CONATIVE, BODY-PART POSSESSOR, INSTRUMENT SUBJECT, INTRANSITIVE INCHOATIVE

To set up the formal context we take the verbs as objects. The properties are the alternations and the relation specifies whether, for a given verb, an alternation is acceptable or not. The resulting formal context is shown in Table 3.

| | Instrument Subject | Middle | Trans. Body-part or refl. | Intransitive Inchoative | Conative |
|----------|--------------------|--------|---------------------------|-------------------------|----------|
| break | x | x | | x | |
| splinter | x | x | | x | |
| smash | x | x | | x | |
| cut | x | x | x | | x |
| scratch | x | x | x | | x |
| scrape | x | x | x | | x |
| hit | x | | x | | x |
| strike | x | | x | | x |
| bash | x | | x | | x |
| touch | | | x | | |
| sting | | | x | | |
| kiss | | | x | | |

Table 3: Formal context for 12 break, cut, hit and touch verbs with Beth Levin's diathesis alternations.

Figure 1 shows the concept lattice obtained when applying formal context analysis to this formal context.

We will explain the structure and labeling of the concept lattice more in detail in Section 5. But Figure 1 already shows that the 12 verbs were divided in 4 classes and

that the classes correspond to Beth Levin’s and the VerbNet classification presented in Section 2.1.

A further interesting topic arising with FCA is identifying dependencies among attributes. Attribute exploring is concerned with the following questions: Does an attribute set imply another one? What is the probability that an attribute set implies another attribute set (given this particular data)? Applied to our problem we will see that FCA defined devices such as attribute association rules help find alternations for French, given the type of verb lexicons mentioned in Section 4.

In this work we show how the underlying idea of the FCA approach, namely identifying classes based on objects sharing a maximal set of attributes can be exploited in several ways for building a classification of French verbs:

1. We compute verb classes using the verbs and properties provided by the French resources described in Section 4

by building FCA lattices.

2. We identify diathesis alternations by looking for the most relevant pairs of attributes

using FCA association rules.

3. We try to match the computed concepts and VerbNet classes by identifying measures for the similarity of their defining sets of objects and attributes

using graph matching algorithms for comparing FCA lattices with other hierarchical structures.

4 The Resources

In this section we briefly describe the three resources we used as input for our classification. For each resource, we give some historical and theoretical background, information about its adequacy for NLP use and general conditions of use. We also show some sample entries.

Information about coverage is summarised in the last part of this section (Subsection 4.4).

4.1 Volem

Linguistically, Volem ([FVSD⁺02]) is closest to Beth Levin’s verb classes. It is a resource built manually within the regional European project VOLEM and building on earlier work by Patrick Saint-Dizier ([SD96], [SD99]). The methodology is also very similar to that of Beth Levin: Patrick Saint-Dizier and his collaborators defined a set of diathesis alternations (called *contexts* in this framework) and assigned them to the verbs, according to whether the verb’s usage in this *context* is acceptable or not. The initial set of *contexts* was enlarged and unified, in order to account for Romance languages other than French.

The resource is available from the authors as an xml file, a sample entry of which is shown below:

```
<TERME>  
<VERBE>supprimer</VERBE>  
<LCS />
```

```

<ROLETHERM>[[inic(agent)],[tid,tiv]]</ROLETHERM>
<ALTERNANCES>
  caus_2np,anti_pr_np,pas_etre_part_np_pp,
  pas_etre_part_np,state_2np_pp,caus_refl_pr_np
</ALTERNANCES>
<ALT-ANCIENNES>12,50,60,162,172</ALT-ANCIENNES>
<WN>8,2</WN>
<EXEMPLE>Ils ont supprimé ce mur</EXEMPLE>
</TERME>

```

For the verb classification, we used the verbs and the context identifiers provided in the ALTERNANCES and ALT-ANCIENNES elements.

4.2 The Ladi Grammar-Lexicon

As mentioned in Section 4, the Grammar-Lexicon, is a resource developed by Maurice Gross and his collaborators at the “Laboratoire d’Automatique Documentaire et Linguistique” from 1968 to 2002 ([Gro75], [BGL76] and [BGL76]). Its initial purpose was to provide a systematic description of the syntactic properties of syntactic functors for French, in particular verbs. The Grammar-Lexicon consists of a set of tables where each table groups together (usages of) predicative items that share some definitional properties. In particular, all predicative items in a given table share one (sometimes two) basic constructions (subcategorisation frames).

| N0 =: Nhum | N0 =: Nnr | N0 =: le fait Qu P | N0 = V1 W | N0 = V2 W | N0 = V2c W | 6 | 10 | N0 V | N0 V N1 | N1 =: Qu P | N1 =: Qu Psubj | N1 = aux V0 W | N1 = de V0 W | N1 = de V2 W | N1 = de V2c W | N1 =: si P ou si P | N1 =: ee(ci+la) |
|------------|-----------|--------------------|-----------|-----------|------------|---|--------------|------|---------|------------|----------------|---------------|--------------|--------------|---------------|--------------------|-----------------|
| - | + | + | - | + | + | | pré supposer | - | + | - | - | + | - | + | - | - | + |
| - | + | + | - | + | + | | supposer | - | + | + | - | + | - | + | - | - | + |
| - | + | + | - | - | + | | supprimer | - | - | + | - | + | - | - | + | - | + |
| - | + | + | + | + | + | | supprimer | - | + | - | - | + | - | + | - | - | + |

Table 4: Sample rows of the Ladi table 10

Table 4 shows some sample lines of the Ladi table 10.

The grammar lexicon has been digitised by the Laboratoire d’Automatique Documentaire et Linguistique (LADL) and is now partially available⁸ under an LGPL-LR licence⁹. Yet its use within natural language processing systems is still hampered both by its non standard encoding and by a structure that is partly implicit and partly underspecified. The SynLex subcategorisation lexicon ([GGPF06]) translates this information into a format more amenable for use by NLP systems and is currently being validated by human annotators within the BDSyn project¹⁰.

In this thesis, we only use the information whether a verb is a member of a table. A possible next step would be to use other properties provided by the tables or the subcategorisation frames of SynLex.

The coverage is summarised in Section 4.4.

4.3 Dicovalence

The syntax lexicon Dicovalence is developed at the University of Leuven (Belgium) by Karel van den Eynde and Piet Mertens, following on initial work from 1986 to 1992

⁸~ 60% of the table/verb information

⁹cf. <http://infolingu.univ-mlv.fr/DonneesLinguistiques/Lexiques-Grammaires/Presentation.html>

¹⁰<http://www.loria.fr/~guillaume/BDSyn/index.html>

within the Proton research project.

This syntax lexicon lists the “valency frame(s)” of French verbs. A valency frame is the set of subcategorized terms (complements and subject), indicating their syntactic function and some of their properties.

In Dicovalence there is an entry for each verb and all the constructions this verb can enter. A sample is shown below:

| | |
|----------|--|
| VAL\$ | supprimer: P0 P1 |
| VTYPES\$ | predicator simple |
| VERBS\$ | SUPPRIMER/supprimer |
| NUM\$ | 80500 |
| EG\$ | nous avons supprimé tous les obstacles à la publication de ce dico |
| TR.DU\$ | afschaffen, opheffen, intrekken, weghalen, weglaten, schrappen, doen verdwijnen |
| TR.EN\$ | suppress |
| P0\$ | que, qui, je, nous, elle, il, ils, on, ça, celui-ci, ceux-ci |
| P1\$ | que, qui, te, vous, la, le, les, se réfl., se réc., en Q, ça, ceci, celui-ci, ceux-ci, l'un l'autre |
| RP\$ | passif être, se passif, se faire passif |

The fields VAL, VTYPES, VERB, NUM, EG and TR are present in all entries. Their significance is as follows:

| | |
|----------|--|
| VAL\$ | predicate and short notation for the syntactic construction |
| VTYPES\$ | predicate type: ordinary, auxiliar, resultativ construction |
| VERBS\$ | verb infinitiv, in all caps and without accents/then ordinary spelling |
| NUM\$ | an identifier for the entry |
| EG\$ | an example |
| TR.DU | dutch translation |

Thus, according to this entry, the verb *supprimer* can occur in a construction P0 supprimer P1, where P0 and P1 can be replaced by one of the entries in fields P0\$ and P1\$ respectively.

For our classification task we used the verbs (as objects) and the constructions in short notation of field VAL (P0 P1 for the sample entry).

The resource is available as open software under the conditions of the LGPL-LR license.

4.4 Coverage Summary

The following table gives an overview of the coverage of the 3 resources used in this thesis for building verb classes:

| | verbs | | attributes |
|--------------------|-------------|--------|------------------|
| | number | number | type |
| Volem | 1635 | 101 | contexts |
| Ladl | 5516 | 57 | tables |
| Dicovalence | 3700 | 312 | valency frames |
| VerbNet | 5000 senses | 392 | syntactic frames |

At first glance, Volem seems the most adequate resource for our classification task because of its linguistic similarity with VerbNet. But, as the table above shows, its coverage is very small compared to that of the other resources and is insufficient for a larger scale classification.

5 Formal Concept Analysis

Formal Concept Analysis (FCA) is a method for deriving conceptual structures out of data, used in applications for data analysis, knowledge representation and information management. This and other related keywords were coined in the early 80s by Rudolf Wille [Wil82], but some fundamental notions appeared independently in the 70s in France [BM70].

In this section we briefly introduce the formal concepts underlying FCA.

5.1 Formal concepts

FCA starts from a *formal context*:

Definition 1 (Formal context) A formal context is a triplet $\mathbb{K} := (\mathcal{O}, \mathcal{A}, R)$, where \mathcal{O} is a set of objects, \mathcal{A} a set of attributes, and R is a binary relation between \mathcal{O} and \mathcal{A} (i.e. $R \subseteq \mathcal{O} \times \mathcal{A}$). $(o, a) \in R$ is read as object o has attribute a .

A formal context is often represented as a “cross table”, as in Table 3, where the verbs are the objects, the attributes the diathesis alternations and a verb “has” a diathesis alternation if it enters the corresponding constructions.

Definition 2 (Derivation operator) For a set $O \subseteq \mathcal{O}$ of objects, let

$$O' := \{a \in \mathcal{A} \mid (o, a) \in R, \text{ for all } o \in O\}$$

Dually, for a set $A \subseteq \mathcal{A}$ of attributes, let

$$A' := \{o \in \mathcal{O} \mid (o, a) \in R, \text{ for all } a \in A\}$$

The operator $'$ in O' and A' is called the derivation operator for the formal context $(\mathcal{O}, \mathcal{A}, R)$.

If O is a set of objects, then O' is a set of attributes to which we can apply the second derivation operator to obtain O'' , a set of objects. Dually, this applies to a set of attributes A .

Definition 3 (Formal concept) Let $\mathbb{K} = (\mathcal{O}, \mathcal{A}, R)$ be a formal context, O' and A' the derivation operator applied to $O \subseteq \mathcal{O}$ and $A \subseteq \mathcal{A}$ respectively.

A formal concept C of \mathbb{K} is a pair (O, A) with $O \subseteq \mathcal{O}$, $A \subseteq \mathcal{A}$ and $O' = A$ and $A' = O$. O is called the extent, and A the intent of C .

For a formal concept (O, A) , O and A are the *maximal* subsets of \mathcal{O} and \mathcal{A} respectively, such that $(o, a) \in R$ for every $o \in O$ and $a \in A$.

The description of a concept (O, A) by extent and intent is redundant, because each of the two components determines the other (since $O' = A$ and $A' = O$), but proved very convenient.

Definition 4 (Subconcept-superconcept relation) Let $\mathbb{K} = (\mathcal{O}, \mathcal{A}, R)$ be a formal context and (O_1, A_1) , (O_2, A_2) be formal concepts of \mathbb{K} . We say that the concept (O_1, A_1) is smaller than (O_2, A_2) :

$$(O_1, A_1) \leq (O_2, A_2) \text{ iff } O_1 \subseteq O_2$$

Which is equivalent to:

$$A_1 \supseteq A_2$$

This relation defines a partial order on the formal concepts.

Definition 5 (Infimum) Let $C_1 = (O_1, A_1)$ and $C_2 = (O_2, A_2)$ be formal concepts of some formal context and $'$ be the derivation operator. The infimum of the formal concepts C_1 and C_2 is defined as:

$$(O_1, A_1) \wedge (O_2, A_2) := (O_1 \cap O_2, (A_1 \cup A_2)'')$$

Remark 1 For every two concepts C_1 and C_2 there exists a unique infimum. It is also called the **the greatest common subconcept** of C_1 and C_2 .

Definition 6 (Supremum) Let $C_1 = (O_1, A_1)$ and $C_2 = (O_2, A_2)$ be formal concepts of some formal context and $'$ be the derivation operator. The supremum of the formal concepts C_1 and C_2 is defined as:

$$(O_1, A_1) \vee (O_2, A_2) := ((O_1 \cap O_2)'', A_1 \cup A_2)$$

Remark 2 For every two concepts C_1 and C_2 there exists a unique supremum. It is also called the **the least common superconcept** of C_1 and C_2 .

The existence of infima and suprema imply that the set of concepts with the partial order given by the relation defined in Def. 4 is a complete lattice:

Lemma 1 (Concept lattice) The set of all formal concepts of a context \mathbb{K} together with the order relation \leq is a complete lattice, i.e. for each subset of concepts there is always a unique greatest common subconcept and a unique least common superconcept.

This lattice is called the concept lattice of \mathbb{K} .

For every object o , the concept $(\{o\}'', \{o\}')$ is the “smallest” concept of the lattice for which the extent contains o . For every attribute a , the concept $(\{a\}', \{a\}''')$ is the “greatest” concept of the lattice for which the intent contains a . Not every concept can be written as $(\{o\}'', \{o\}')$ for an object o or $(\{a\}', \{a\}''')$ for an attribute a , but if there are objects and/or attributes with this property the concept typically is labeled with them:

Definition 7 (Concept labeling) Let $\mathbb{K} = (\mathcal{O}, \mathcal{A}, R)$ be a formal context, $o \in \mathcal{O}$ an object, $a \in \mathcal{A}$ an attribute. In the concept lattice of \mathbb{K} we label with:

$$\begin{aligned} o \text{ the concept } (\{o\}'', \{o\}')$$

$$\rightarrow \text{ object labels,}$$

$$a \text{ the concept } (\{a\}', \{a\}''')$$

$$\rightarrow \text{ attribute labels.}$$

Definition 8 (Reduced extent and intent) Let $\mathbb{K} = (\mathcal{O}, \mathcal{A}, R)$ be a formal context, C a concept of \mathbb{K} .

The set of object labels of C is called the **reduced extent** of C .

The set of attribute labels of C is called the **reduced intent** of C .

Figure 1 shows the concept lattice of the context in Table 3 by a *line diagram* or *Hasse diagram*. The diagram has been generated by the ConExp software¹¹.

In a line diagram each vertex represents a formal concept. A concept C_1 is a subconcept of a concept C_2 iff there is a descending path from vertex C_2 to vertex C_1 . Each vertex (formal concept) C is labeled with its object and attribute labels (if there are any). If a vertex C has an object label o , it is the **smallest** concept with o in its extent; dually, if the vertex C is labeled with an attribute a , it is the **largest** concept with a in its intent.

The ConExp lattice drawer represents concepts as circles. It will fill the upper half of a circle representing a concept if this concept has an attribute label. The lower half of the vertex is filled when the concept is labeled with an object. Attribute labels are mouse coloured and attached above and object labels are white and attached below the concept vertex.

¹¹<http://conexp.sourceforge.net/projects/conexp>

In Figure 1 there are two smallest concepts, labeled with objects and attributes as follows:

| concept id | object labels | attribute labels |
|------------|------------------------|-------------------------|
| C_1 | cut, scrape, scratch | – |
| C_2 | break, splinter, smash | INTRANSITIVE INCHOATIVE |

Thus, the reduced extent of C_1 are the verbs *cut*, *scrape*, *scratch* and of C_2 *break*, *splinter*, *smash*. C_1 has no reduced intent, the reduced intent of C_2 is INTRANSITIVE INCHOATIVE. Because C_1 and C_2 are smallest concepts, their reduced extent is also their extent. To know the intent of eg. concept C_1 we have to follow all ascending paths to the top concept and collect the attribute labels: CONATIVE, TRANS BODY PART OR REFLEXIVE OBJECT, MIDDLE, INSTRUMENT SUBJECT.

We will see that for our use case, the reduced extents yield the verb classes. By contrast, we found that a reduced intent of more than one element typically indicates irregularities in the data.

5.2 Association Rules

Dependencies between attributes can be described by *implications* and *association rules*.

Definition 9 (Itemset) Let $\mathbb{K} = (\mathcal{O}, \mathcal{A}, R)$ be a formal context. A set of attributes $A \subset \mathcal{A}$ is called an *itemset* of \mathbb{K} .

Definition 10 (Support) Let A be an itemset, A' the derivation operator:

The *support* of A is:

$$\text{support}(A) := |A'|$$

that is, the number of objects which have this itemset.

Definition 11 (Frequent itemsets) An itemset is said to be *frequent* if its support is above a user defined threshold.

In the example in Figure 1 we have:

$$\begin{aligned} \text{support}(\text{MIDDLE}) &= 6, \text{ and} \\ \text{support}(\text{INSTRUMENT SUBJECT}) &= 9. \end{aligned}$$

Definition 12 (Association rule, confidence) Let A_1, A_2 be itemsets of a formal context \mathbb{K} .

An expression of the form $A_1 \rightarrow A_2$ is called an *association rule*

The *confidence* of the association rule $A \rightarrow B$ is defined as:

$$\text{confidence}(A_1 \rightarrow A_2) := \frac{\text{support}(A_1 \cup A_2)}{\text{support}(A_1)}$$

For instance, considering the association rule $r = \text{MIDDLE} \rightarrow \text{TRANS BODY-PART OR REFL}$ we find:

$$\text{confidence}(r) = \frac{\text{support}(\text{MIDDLE}, \text{TRANS BODY-PART OR REFL})}{\text{support}(\text{MIDDLE})} = \frac{3}{9} = 33\%$$

Let us now consider the confidence of the association rule $r = \text{MIDDLE} \rightarrow \text{INSTRUMENT SUBJECT}$:

$$\text{confidence}(r) = \frac{\text{support}(\text{MIDDLE, INSTRUMENT SUBJECT})}{\text{support}(\text{MIDDLE})} = \frac{6}{6} = 1$$

A rule of confidence 1 is called **valid**, it is an implication and valid globally in the context:

Definition 13 (Implication) Let $\mathbb{K} = (\mathcal{O}, \mathcal{A}, R)$ be a formal context, $A_1, A_2 \subset \mathcal{A}$. We say the **implication** $A_1 \rightarrow A_2$ holds in context \mathbb{K} if $A_1' \subseteq A_2'$, where $'$ is the derivation operator introduced in Definition 2.

In other words, all objects which share attributes A_1 also have attributes A_2 .

Thus, we already saw that the implication $\text{MIDDLE} \rightarrow \text{INSTRUMENT SUBJECT}$ holds for the context corresponding to the concept lattice in Figure 1.

The confidence can be seen as a conditional probability: The probability that an object of this context has all the attributes in A_2 , given that it has all the attributes in A_1 . However, the confidence of a rule only is meaningful in combination with its support¹²: A rule may have a confidence of one, but this may not be of much interest if the support is only 1 or 2, i.e. there are very few objects with this attribute. Nonetheless, the sheer existence of the rule may be of some importance by itself.

Definition 14 (equivalence) Let A_1, A_2 be itemsets of a formal context \mathbb{K} .

A_1, A_2 are equivalent, $A_1 : \leftrightarrow A_2$, iff $A_1 \rightarrow A_2$ and $A_2 \rightarrow A_1$ are implications.

Remark 3 Let a_1 and a_2 be distinct attributes in a reduced intent. Then $\{a_1\}$ and $\{a_2\}$ are equivalent: $\{a_1\} \leftrightarrow \{a_2\}$.

Equivalent attributes can be used to simplify a context: they can be replaced by a single attribute, the name of which is (for example) the list of the equivalent attribute's names.

Application to verb classes We will apply FCA as indicated in the example: The raw data are the syntax lexicons presented in Section 4. In all the experiments the verbs will be considered as objects. The attributes will be syntactic properties provided by these resources. First, we build classes and second, we investigate the distribution of association rules of length 2 (between 2 attributes) to hopefully obtain relevant alternation candidates.

6 Extraction of Verb Class Information from Syntactic Lexicons

In this section we apply FCA to the French syntax lexicons described in Section 4 for building verb classes (Section 6.1) and mining for alternation candidates (Section 6.2).

In Subsection 6.3 we discuss the outstanding issues arising with these enterprises and the direction for further investigations with regard to evaluation that are suggested by these issues.

¹² $\text{support}(A_1 \rightarrow A_2) := \text{support}(A_1 \cup A_2)$

6.1 Building Verb Classes

We experimented with a straight forward way of building classes from the resources described previously (Volem 4.1, the Ladl Grammar-Lexicon 4.2 and Dicovalence 4.3), following an FCA based technique described in [CHS05].

- The verbs of the different resources are the *objects* of the formal context.
- The *attributes* of the formal context are the syntactic properties assigned to the verbs by these resources. That is, *contexts* for Volem 4.1, a specific *table* for the Ladl Grammar-Lexicon 4.2 and a *valency frame* for Dicovalence 4.3¹³.

For illustration, consider the sample entries shown in Section 4. They display the following object – attribute assignments:

Volem: the verb is **supprimer**, assigned *attributes* are: **caus_2np, anti_pr_np, pas_etre_part_np_pp, pas_etre_part_np, state_2np_pp, caus_refl_pr_np, 12, 50, 60, 162, 172**

Ladl: the verbs are **présupposer, supposer** and **supprimer**, the assigned *attribute* is **10** (the table number).

Dicovalence: the verb is **supprimer** and the assigned *attribute* is **P0 P1** (the valency frame in short notation, as shown in the field VAL\$).

Table 16, 17 and 18 in the Appendix illustrate the way we build formal contexts from the Volem, Ladl and Dicovalence resources respectively.

From these formal contexts we build a concept lattice for each resource. We then use the concept lattice to extract verb classes along the lines of the following considerations, which once more highlight the particular suitability of FCA for this task.

A formal concept of the lattice associates a maximal set of objects with a maximal set of attributes. That is, a tuple (O, A) of objects and attributes is a formal concept iff A consists of all and only the attributes true of all the objects in O and conversely, there is no object not in O of which all the attributes in A are true. In other words, FCA directly gives us the “best” correspondence between “maximal” verb and attribute sets. For instance, given the context in Table 3, FCA yields the concepts in Table 5.

We see that there are no two concepts with identical alternation set but distinct verb set or conversely, identical verb set but distinct alternation set. That is, FCA allows us to uniquely identify the verb classes defined by alternations.

Nonetheless, the classes thus defined are not entirely appropriate as they fail to bring out clearly groups of verbs whose syntactic behaviour differs. For instance, we can see from the above list of concepts that *Break* verbs accept the INSTRUMENT SUBJECT, the MIDDLE and the INTRANSITIVE INCHOATIVE alternation whereas *Cut* verbs only accept INSTRUMENT SUBJECT and MIDDLE. Hence we would like *Cut* and *Break* verbs to form distinct classes namely, a class including only the *Cut* verbs and accepting INSTRUMENT SUBJECT and MIDDLE and another class including only the *Break* verbs which would additionally accept the INTRANSITIVE INCHOATIVE alternation. In fact, FCA gives us the possibility to automatically identify such classes by considering the reduced extents rather than the full ones. Thus if instead of listing the full extent of a concept we consider its reduced extent, we get the list of concepts in Table 6¹⁴.

¹³We sometimes denote syntactic properties by *attributes* or *frames*

¹⁴Note that the concept lattice in Figure 1 shows only the reduced extents and intents, whereas Table 5 lists the full extents and intents

| class name | verbs | properties |
|------------|--|---|
| C_1 | cut, scrape, scratch | TRANS BODY PART OR REFLEXIVE OBJECT, CONATIVE, MIDDLE, INSTRUMENT SUBJECT |
| C_2 | break, splinter, smash | INTRANSITIVE INCHOATIVE, MIDDLE, INSTRUMENT SUBJECT |
| C_3 | hit, strike, bash, cut, scrape, scratch | TRANS BODY PART OR REFLEXIVE OBJECT CONATIVE, INSTRUMENT SUBJECT |
| C_4 | cut, scratch, scrape, break, splinter, smash | MIDDLE, INSTRUMENT SUBJECT |
| C_5 | touch, sting, kiss, hit, strike, bash, cut, scrape, scratch | TRANS BODY PART OR REFLEXIVE OBJECT |
| C_6 | hit, strike, bash, cut, scrape, scratch, break, smash, splinter | INSTRUMENT SUBJECT |

Table 5: Classes obtained from the concept lattice in Fig. 1, full extent.

| class name | verbs | properties |
|------------|------------------------|---|
| C_1 | cut, scrape, scratch | TRANS BODY PART OR REFLEXIVE OBJECT, CONATIVE, MIDDLE, INSTRUMENT SUBJECT |
| C_2 | break, splinter, smash | INTRANSITIVE INCHOATIVE, MIDDLE, INSTRUMENT SUBJECT |
| C_3 | hit, strike, bash | TRANS BODY PART OR REFLEXIVE OBJECT CONATIVE, INSTRUMENT SUBJECT |
| C_4 | – | MIDDLE, INSTRUMENT SUBJECT |
| C_5 | touch, sting, kiss | TRANS BODY PART OR REFLEXIVE OBJECT |
| C_6 | – | INSTRUMENT SUBJECT |

Table 6: Classes obtained from the concept lattice in Fig. 1, reduced extent.

By extracting only the concepts with a non empty reduced extent, we obtain exactly the groupings we are looking for, namely the maximal sets of verbs which accept a given set of alternations.

The concept lattices can be transformed into hierarchical classifications by removing the bottom element. The hierarchy obtained from the concept lattice in the running example is shown in Figure 2(a) and Figure 2(b) shows the hierarchy when concepts with empty reduced extent are pruned.

We give some samples of the classes finally obtained in Table 15 (in the Appendix).

Table 7 gives some features of the classifications we obtained for Volem, the Grammar-Lexicon and Dicovalence. To allow for an estimative comparison with VerbNet, we also show class, member and frame counts for VerbNet. Note however, that it is not entirely clear what a syntactic frame for VerbNet is: Depending on whether we take into account selectional and/or syntactic restrictions we obtain different frames and frame counts. For these figures we only considered the frame name and its syntax description (as given in Table 2).



(a) The classes are computed using the full or reduced extent. The class labels are those assigned in Table 5. (b) Classes are computed using reduced extents, classes with empty reduced extent were discarded.

Figure 2: Class hierarchies obtained from the concept lattice in Fig. 1. Classes are labeled as in Table 5

| Resource | verb nbr. | attribute nbr. | class nbr. | classes with 1 frame | classes with 1 verb | largest class | largest frame set | average depth | longest path | average class size | median class size | average frame nbr. | median frame nbr. |
|-----------------|-----------|----------------|------------|----------------------|---------------------|---------------|-------------------|---------------|--------------|--------------------|-------------------|--------------------|-------------------|
| Volem | 1635 | 101 | 936 | 2 | 742 | 68 | 26 | 2.13 | 8 | 1.75 | 1 | 12.7 | 13 |
| Grammar-Lexicon | 5516 | 57 | 1778 | 56 | 1415 | 253 | 18 | 2.11 | 4 | 3.1 | 1 | 3.75 | 3 |
| Dicovalence | 3737 | 312 | 1008 | 3 | 823 | 908 | 38 | 2.4 | 6 | 3.7 | 1 | 8.1 | 7 |
| VerbNet | 3626 | 392 | 430 | 34 | 53 | 383 | 26 | 1.57 | 4 | 12.63 | 6 | 3.7 | 3 |

Table 7: Some features of the verb class hierarchies built from Volem 4.1, the Grammar-Lexicon 4.2 and Dicovalence 4.3 as compared to VerbNet.

Class hierarchy We see from Table 7 that the VerbNet hierarchy is very flat. The hierarchies of the acquired classifications are also rather flat, but deeper than the VerbNet hierarchy.

In the following we examine how the verbs are distributed among classes characterised by their size (the number of verbs they contain – Figure 3) and their frame set cardinality (Figure 4).

Figure 3 shows the distribution of the verbs in the classifications obtained from Volem, Ladd and Dicovalence on one hand and VerbNet on the other in classes of a specific size (number of verbs). The x-axis displays the class size and the y-axis the proportion of verbs (in %) assigned to a class of the size given by the x-axis. Similarly as was remarked previously with regard to Table 7, the delimitation of the exact number of syntactic frames for VerbNet is difficult to determine.

First, we observe that generally speaking, the verbs are more uniformly distributed in the VerbNet classification than in any of the computed classifications. The computed classifications have a very large number of classes and a large proportion of verbs occur in small classes. In addition, the Dicovalence verb distribution is also skewed towards the large classes.

Small classes. We see that in the computed classifications most verbs are assigned to small classes (< 5 members). In particular, Table 7 shows that most of these consist

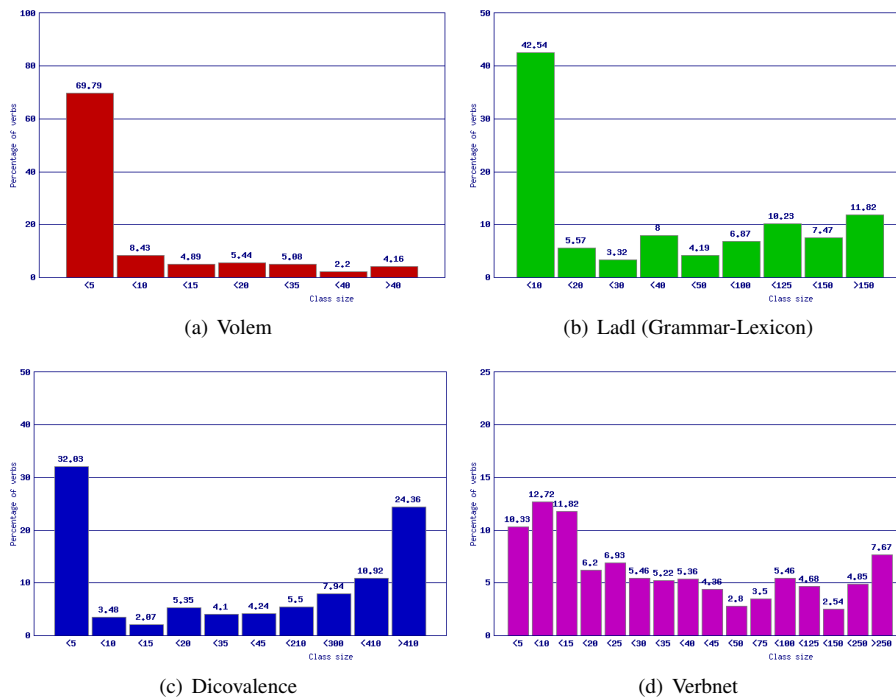


Figure 3: Distribution of verbs in classes wrt. the size of the classes. The figures (a), (b) and (c) show the verb distribution for the computed classifications for Volem (4.1), the Grammar-Lexicon (aka Ladl, 4.2) and Dicovalence (4.3) respectively. They show the percentage of verbs pertaining to a class of size given by the x-axis. To allow a comparison with VerbNet we also give VerbNet counts (2.1.2) in (d). The x-axis gives the class size. The y-axis shows the percentage of verbs which are in a class of the size given by the x-axis.

of only one verb. It is not clear at this stage whether this is due to our use of the information or to the quality and quantity of the data: some classes may only contain one verb simply because other relevant verbs are not present in the resources used.

Large classes. In the Dicovalence classification an important proportion of verbs were assigned to large classes, in particular to one very large class of 908 members. The corresponding intent of this class consists of the single valency frame P0 P1 – the basic transitive construction, which naturally is shared by a large number of verbs. This reinforces the approach to diathesis alternations discussed in Section 6.2: Diathesis alternations mostly involve a common basic construction (like TRANSITIVE or INTRANSITIVE) which is varied by the second involved construction. It seems more meaningful to base the classification on whether a verb accepts a pair of such constructions rather than one very common construction.

On the other hand, we see that, surprisingly in Dicovalence, 908 verbs only accept this single very common construction. Going more in detail, we found that in fact this construction is further differentiated by the possible realisations of P0 and P1. This shows that in this case our choice of attributes was obviously too rough and therefore inadequate.

Large number of classes. Table 7 shows that the number of classes in the computed classifications is very large. This is a known and recurrent issue with FCA. Typically, it is addressed by reducing the complexity of the problem eg. by means of the FCA technique of *Iceberg Lattices* (cf. [STB⁺02] and [BTP⁺02]). This approach is based on *frequent itemsets*. The constructed lattices consist only of the top most concepts, whose intent cardinality is above a minimum support threshold: These are the concepts which provide the most global structuring of the domain.

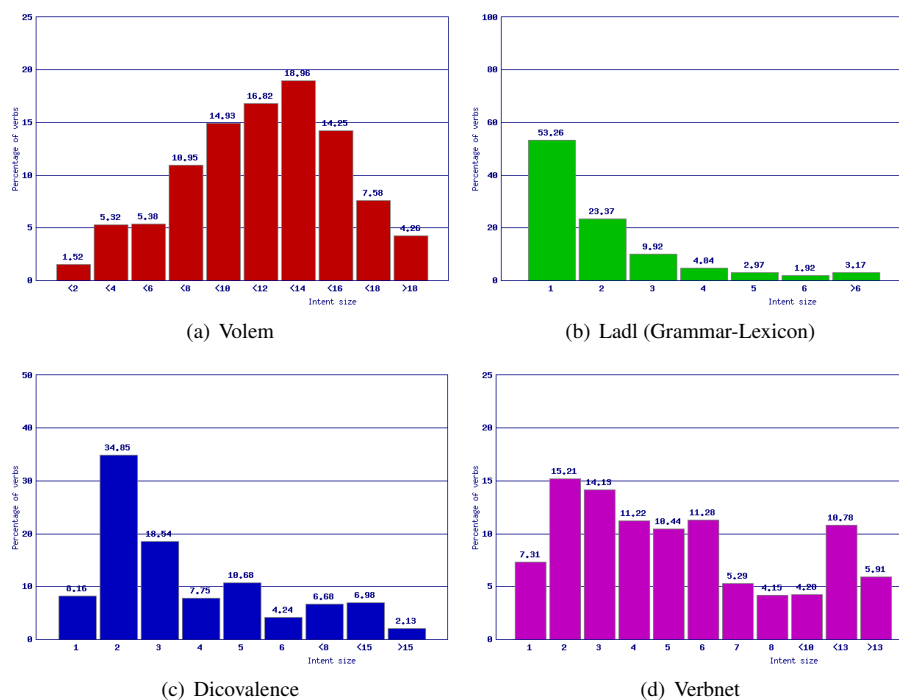


Figure 4: Distribution of verbs in classes wrt. the size of the class frame set. The figures (a), (b) and (c) show the verb distribution for the computed classifications for Volem (Sec. 4.1), the Grammar-Lexicon (aka Ladl, Sec. 4.2) and Dicovalence (Sec. 4.3) respectively. We count the percentage of verbs which are in classes, the frame set of which has a specific size, given on the x-axis. To estimate the analogy to VerbNet, we also give counts for VerbNet (Sec. 2.1.2) in (d). The y-axis shows the percentage of verbs which are in a class which has a frame set of the size given by the x-axis.

Figure 4 shows the repartition of verbs in classes characterised by the size of their intent. The x-axis gives the cardinality of the classes' frame sets and the y-axis the proportion of verbs (in %) assigned to a class with a frame set size given by the x-axis. Again, we observe that the distribution for VerbNet is more uniform than that of the computed classifications. VerbNet has a rather balanced verb repartition with a slight bias on classes with 2 frames, whereas the Volem classification assigns most verbs to classes with 13-14 constructions, the Ladl classification to classes with 1 table, and most Dicovalence verbs are in classes with 2 valency frames. An interesting issue is the proportion of verbs in classes with a single frame as intent – in a way these shouldn't be considered as proper classes because a diathesis alternation cannot consist of only one frame. We see that only Ladl displays an unusual proportion of verbs in such classes.

There are several possible reasons: First, only 60% of the verb-table information is available. Second, variations similar to diathesis alternations are also expressed inside a table, which we didn't take into account. Finally, some of the Ladr tables correspond to two basic constructions.

Conclusion. These simple observations derived from descriptive statistics show that the acquired classifications are very different from the targeted ones. Nevertheless they represent a starting point, raise some interesting issues and give clues to directions for further work.

6.2 Mining Diathesis Alternations for French

In this section we will use the FCA notion of *frequent itemsets* (cf. Section 5.2) to find candidates for diathesis alternations for French. These alternations were unified and validated manually by Claire Gardent and Fabienne Venant and aligned as far as possible with the VerbNet alternations. To find relevant diathesis alternations for French is an interesting task in itself¹⁵, but a further objective is to build classes by again considering the verbs as objects, and the observed alternations as attributes. In turn, these classes could then provide further insights with regard to the validity of the alternation candidates.

The methodology is the following:

1. We consider the resources **Volem** (Section 4.1) and the **Grammar Lexicon** (Section 4.2) separately¹⁶.
2. As in Section 6.1, the verbs are the objects and the attributes are *contexts* in the case of **Volem** and *table names* in the case of the **Grammar Lexicon**. We are looking for association rules $a \rightarrow b$, where a and b are attributes such that $support(a)$ and $confidence(a \rightarrow b)$ are high.
3. We produced a list of attributes, sorted by descending support. For each of these attributes we list the rules $a \rightarrow b$ by descending confidence, down to a confidence of 20%. To facilitate the decision whether an association rule is effectively a diathesis alternation, we display the verbs having these attributes and available examples.
4. The experts choose the linguistically most coherent alternation candidates from the list produced in the preceding step.
5. We create a context with the verbs as objects and the selected diathesis alternations as attributes, build the concept lattice and extract classes, as in Section 6.1.

At the time of writing, we implemented stages 1 to 3 and are in the course of processing stage 4.

Table 8 shows an excerpt of the list of alternation candidates for **Volem**¹⁷. In this table we inspect the pairs of attributes a, b , where $a = anti_pr_np$. *anti_pr_np* has a support of 1162, that is there are 1162 (of 1635) verbs which have this attribute. The next three rows in the table (set in bold face) show an example for this attribute and a

¹⁵Diathesis alternations are interesting from a linguistic point of view because they represent a linguistic primitive of a language. On the other hand, as we saw earlier, they are cornerstones in building syntactic-semantic verb classes.

¹⁶In future work we plan to also use this approach with **Dicovalence** (Section 4.3).

¹⁷The complete tables are available at http://verbnet-fr.gforge.inria.fr/alternation_candidates_volem_old_new_vn.html.

frame name and syntax description determined by linguistic experts as close as possible to VerbNet’s terminology.

The next row shows that the rule $a \rightarrow b$ with $b = pas_etre_part_np_pp$ has a support of 1125 and a (high) confidence of 0.9673. *pas_etre_part_np_pp* designates a prepositional passive construction, in VerbNet terminology *PP Passive*, and the syntax description is *Patient be V by Agent*.

| Premise attribute: <i>anti_pr_np</i> (1162 verbs) | | |
|---|------------|-----------------------------------|
| Example(s) | | La porte s’ouvrit |
| VerbNet frame name | | Intransitive Inchoative-Reflexive |
| VerbNet Syntax description(s) | | Patient {refl} V |
| Rule $a \rightarrow b$ | confidence | support(a,b) |
| <i>pas_etre_part_np_pp</i> | | |
| PP Passive | | |
| Patient be V by Agent | 0.9673 | <u>1124</u> |
| La porte est ouverte à Marie par Jean | | |
| <i>caus_2np</i> | | |
| Basic Transitive | | |
| Agent V Patient | 0.9475 | <u>1101</u> |
| Jean ouvre la porte | | |

Table 8: Alternation candidates for **Volem**^a.

^aThe sample constructions have been changed for illustration purposes.

For the sake of completeness we also show some alternation candidates for the **Grammar Lexicon** in Table 9¹⁸.

| Premise attribute: <i>35R</i> , subclass of <i>Basic PP</i> , basic (390 verbs) | | |
|---|------------|------------------------|
| Example(s) | | Paul compte sans Marie |
| VerbNet frame name | | PP |
| VerbNet Syntax description(s) | | Theta1 V Theta2 |
| Rule $a \rightarrow b$ | confidence | support(a,b) |
| Basic transitive | 0.45641 | <u>178</u> |
| Basic intransitive | 0.34615 | <u>135</u> |
| 32R3 | | |
| Basic transitive constrained-OBJECT | | |
| Agent V Patient < +constrained > | 0.27436 | <u>107</u> |
| Paul parle l’anglais | | |
| Basic ditransitive | 0.25128 | 98 |
| 31H | | |
| Basic intransitive human-SUBJECT | | |
| Theta1 < +human > V | 0.21795 | <u>85</u> |
| Paul boîte | | |
| 31R | | |
| Basic intransitive | 0.20256 | <u>79</u> |
| Theta1 V | | |
| Le blé pousse | | |

Table 9: Alternation candidates for **Ladl**.

6.3 Discussion

The conducted experiments (described in Section 6.1 and 6.2) raised 3 outstanding issues, which will be discussed in this section:

1. Evaluation
2. Improving initial resources

¹⁸The complete tables are available at <http://verbnet-fr.gforge.inria.fr/alternation.candidates.ladl.vn.html>.

3. Fusioning resources

Item 1 is of most consequence and will be considered more in detail in the next subsection. We now briefly address the other two topics.

Improving initial resources. As a byproduct of the class building experiments presented in Section 6.1 we identified various irregularities and inconsistencies (for instance frame co-occurrences, typos, etc.), which could be corrected. A more systematic approach to this kind of error mining would be desirable.

Fusioning resources. Put together the three French syntax resources contain 5770 distinct verbs, 1437 verbs are shared by all three resources, 74 are unique to Volem, 1902 to Ladl and 189 to Dicovalence. Considering these counts, it is tempting to merge the resources. This could be done in two ways:

1. By building the context from the union of verbs as objects and the union of syntactic properties as attributes.
2. By building classifications first and trying to merge them afterwards.

wrt. 1 Considering the results of our class building experiments from Section 6.1, we expect FCA to produce classes of unsatisfactory quality and quantity due to the low quality of the original resources. Furthermore, the complexity and size of the concept lattice increase exponentially with the size of the context. Nevertheless, it is feasible, and, as seen in Section 6.1, we might come across interesting implications. The result should be further improved if the 3 resources can be aligned, that is, if the syntactic attributes used by each of the resources are redefined using a common vocabulary. Such validation and normalisation work is in progress (cf. Sec. 6.2).

wrt. 2 This problem is related to ontology merging, an active and productive domain of research. Methods and techniques developed there could induce resolution methods here. Again, the unification of the syntactic properties could help defining similarity measures for classes which in turn can help to match classes.

6.3.1 Evaluation

After computing a classification we face the problem of evaluating it: does it give correct verb classes? The classifications need to be compared at two levels: class-wise and as hierarchical structures.

The following evaluation schemes are conceivable:

Gold Standard. Comparing with regard to a Gold Standard would be a natural way to evaluate our verb classes. However, a suitable Gold Standard is currently inexistent. Here techniques for ontology learning could be used to help create the required reference (See eg. [CHS05], [SWG05]).

If a Gold Standard were available, we could follow the evaluation scheme proposed in [CHS05] (and in [MS02] and [DS06])¹⁹ to judge both the similarity of classes and of the hierarchies.

English VerbNet. A second strategy would be to compare the created lattice to the English VerbNet. Arguments in favour of this approach are that the classes of both resources are determined by similar items: verbs and syntactic constructions. As both English VerbNet classes and the classes we aim at should share semantic primitives we

¹⁹The method has been implemented in the Abraxas project, <http://nlp.shef.ac.uk/abraxas/onteval.html>

presuppose a certain similarity. However, the classes cannot be expected to be congruent and we would have to account for an unknown tolerance range in our calculations. Nevertheless, the verbs and syntactic constructions as class defining devices, provide us with means of relating classes, namely by comparing both their member and frame sets:

1. Verbs can be related by translation. However, translation without careful disambiguation introduces a lot of noise into the data (as the translation Table 12 and a look-up in a dictionary easily show).
2. Linguist experts can assign similarity scores to the syntactic properties of the French resources on one hand and VerbNet's frames on the other.

Thus, despite its shortcomings, this method permits identifying French classes which may either have or lack an English counterpart. Once a counterpart is identified, this can help enlarging existing classes by comparing the involved frames and verbs. If for example a frame is present in a VerbNet class, but not in its French counterpart, we could gain further verb class information by exploring the reasons for this discrepancy. It may be due to

1. a shortcoming in the initial resource, in this case the resource could be improved and the French class would be augmented by the missing frame
2. different syntactic behaviour of otherwise semantically similar French and English classes
3. an incorrect alignment of the two classes by the evaluation method, in this case the evaluation method would have to be adjusted.

Similarly, investigating the differences of the verb members of aligned French and English classes would permit analysing the difference in meaning between French and English verbs and testing the hypothesis that classes share commonalities across languages.

Last but not least, the lack of an English counterpart for a French concept also reveals a relevant linguistic information.

As for **comparing hierarchies**, considering the unidentified structural differences between the two hierarchies, a comparison on the structural level does not seem meaningful for this scheme and at this stage.

Semantic coherence. The most distinctive feature of the classification we would like to achieve is the semantic coherence of its classes. However, to our knowledge there is no automatic means to assess the semantic coherence of a verb class. In particular, the semantic coherence of VerbNet was evaluated manually. A feasible subtask would be an a posteriori evaluation by an expert of a small representative part of the concept lattice.

6.3.2 Conclusion and Direction for further Work

The discussion of the above issues shows that an adequate evaluation methodology for acquired verb class information is essential for all of the tasks described above but also too hard to be solved completely and satisfactorily within the scope of this work. In the following section we set out to explore possible evaluation methods by testing them on a small hand-built example.

7 Possible Evaluation Methodology

In the discussion in Section 6.3 we saw that the most feasible automated evaluation method would be to compare the acquired verb classes with regard to the English VerbNet. To address this issue we identified two techniques for relating concepts to VerbNet classes:

1. using ontology based concept similarity as in [For06].
2. using relational concept analysis as in [HNHRV07],

We tested both techniques on a small hand-built example described in Section 7.2. The outcome lead us to only take into consideration the former method. We briefly motivate this choice in the following.

Method 1. The similarity of an English and a French class is computed as a weighted combination of the similarity of the verb set and the frame set.

Method 2. Relational Concept Analysis (RCA) is a framework which permits enhancing FCA with relational information between objects. Applied to our case, we would relate French verbs to English VerbNet classes by translation.

Method 1 takes both verb and frame similarity into account whereas Method 2 relies entirely on the relation between verbs. Both, the verb similarity in Method 1 and the relation in Method 2 are based on French-English translations. However, Method 1 allows for additional tuning and weighting at different levels: The similarity between two verbs is $\in [0 \dots 1]$ (as opposed to 0 or 1 for Method 2) and can in addition be balanced by the frame similarity.

Our experiments showed that the similarity scores derived from translation are very inaccurate (cf. 12). This, as suggested by the previous observations, had a strong bias on the performance of Method 2, whereas its effect could be balanced in the case of Method 1. As a result, Method 2 gave too unsatisfactory results to be presented in this work.

In the remaining of this section we will therefore focus on Method 1. We will present its theoretical motivation (Section 7.1) and show its operation on a running example (Section 7.2 and 7.3).

7.1 Relating concepts using Ontology Based Concept Similarity

In this section we discuss the Method 1, introduced above.

We adapt a technique described in [For06] and [For08] where the author introduces a notion of similarity between formal concepts of a concept lattice and concepts in a domain ontology. In the following we will first describe this approach before showing how it can be adapted for our purpose.

In this framework, the attributes of a formal concept are used to map it to an ontology concept in the domain ontology²⁰. Attributes are matched to the ontology concept labels and the formal objects correspond to the instances of an ontology concept. The author assumes that ontology concept instances and formal concept objects are represented by elements of the same set, whereas ontology concept labels and formal concept attributes need to be related. This matching is implemented by a similarity

²⁰here the term *concept* is used in two different contexts: in the ontology domain and in the FCA domain. The terms are somewhat related but not at all identical – the same holds for the term *context*

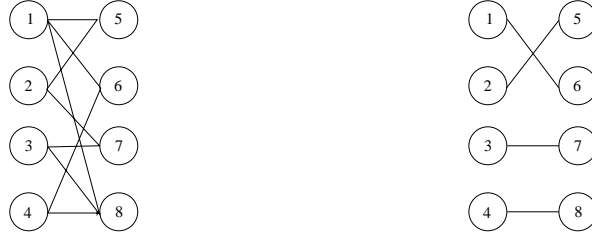
function asd (*axiomatic similarity degree*), on \mathcal{A} , the set of formal attributes and \mathcal{C}_θ , the set of ontology concept labels:

$$asd : \mathcal{A} \times \mathcal{C}_\theta \longrightarrow [0 \dots 1]$$

The similarity function asd is established by domain experts, who assign a similarity score to pairs of ontology concepts and formal attributes.

The author has now the prerequisites to define a similarity function between a formal concept and an ontology concept. The formal concept is given by (E_1, I_1) , its extent E_1 , a set of formal objects, and its intent I_1 , a set of formal attributes. The ontology concept is represented by the pair (E_2, I_2) where E_2 is the set of instances, a set of objects, and I_2 one or more labels of the ontology concept.

The similarity function is inspired by the *maximum weighted matching* problem in bipartite graphs (see [Gal86]). In its simplest form the matching problem may be stated as follows: we are given two sets of vertexes, with arbitrary edges between vertexes from one set to those of the other (Fig. 5). We are interested in particular edge subsets called *matchings*, namely those with edges that have no end point in common. The matching problem consists in finding a matching with maximal cardinality. In the case of the *maximum weighted matching* problem, edges are weighted and we must find a matching with maximal total weight²¹.



(a) Sample input to the matching problem: a graph consisting of 2 sets of vertexes with connecting edges.

(b) A matching of the graph in figure (a). No two edges have an end point in common, there is at most one edge leaving from any vertex.

Figure 5: Illustration of the matching problem in bipartite graphs. Fig. (a) shows a sample input graph and Figure 5(b) a matching. A maximum matching is a matching with the highest number of edges.

Definition 15 (Matchings) Let $n = |I_1|$, $m = |I_2|$ and $n \leq m$.

Then $\mathcal{P}(I_1, I_2)$, the set of all matchings between I_1 and I_2 is defined as follows:

$$\mathcal{P}(I_1, I_2) := \left\{ \{(a_1, b_1) \cdots (a_n, b_n)\} \mid a_h \in I_1, b_h \in I_2, \forall h = 1, \dots, n \right. \\ \left. \text{and } a_h \neq a_k, b_h \neq b_l \forall k, l \neq h \right\} \quad (1)$$

Reformulated in terms of the *matching problem in bipartite graphs*, we consider the graph G with edges (i_1, i_2) for any $i_1 \in I_1$ and $i_2 \in I_2$. Then $\mathcal{P}(I_1, I_2)$ is the set of matchings of G .

²¹ which must not necessarily be unique

Definition 16 (Similarity of a formal concept and an ontology concept) *The similarity of a formal concept (E_1, I_1) and an ontology concept (E_2, I_2) is defined as:*

$$\begin{aligned} \text{Sim}((E_1, I_1), (E_2, I_2)) := & \frac{|E_1 \cap E_2|}{\max(|E_1|, |E_2|)} * w \\ & + \left[\frac{1}{\max(|I_1|, |I_2|)} \max_{P \in \mathcal{P}(I_1, I_2)} \left(\sum_{(a,b) \in P} \text{asd}(a, b) \right) \right] * (1 - w) \end{aligned} \quad (2)$$

where \mathcal{P} is defined as in Definition 15 and $\text{asd}(a, b)$ is the similarity degree as defined by domain experts.

$0 \leq w \leq 1$ is a weight, which can be used for additional tuning.

Finding a matching with a maximum sum in Definition 16 is a reformulation of the *Maximum Weighted Matching in Bipartite Graphs* problem, which is solvable in polynomial time ([Gal86]). It is equivalent to the *Assignment Problem*, which is one of the fundamental combinatorial optimisation problems.

$\text{asd}(a, b)$, the axiomatic similarity degree must not always be defined by domain experts. In a subsequent article, the author uses an Information-Theoretic similarity measure based on the WordNet taxonomy ([Res95] and [Lin98]).

Because of the duality of extent and intent of formal concepts, this similarity measure can easily be adapted to other use cases, as we will see in the following.

In our example the attribute sets are identical. But even if they were not we could still use this similarity definition, provided we could establish a similarity function between the formal attributes on one side and the VerbNet frames on the other.

In our setting the objects (verbs) are equally not identical, only related. We can account for this by introducing an analogous simple similarity function on pairs of French and English verbs.

Definition 17 (Similarity function on French and English verbs) *A very simple similarity function for French and English verbs is:*

$$\text{asd}_{vn}(v_{en}, v_{fr}) := \begin{cases} 1, & \text{if } v_{fr} \text{ is one of the translations of } v_{en} \\ 0, & \text{else.} \end{cases}$$

Finally, we can adapt Definition 16 to our example:

Definition 18 (Similarity of a formal concept and a VerbNet class) *The similarity of a formal concept (E_1, I_1) of French verbs and a VerbNet class (E_2, I_2) is defined as follows. Let*

E_1 be the extent of the formal concept, a set of French verbs.

I_1 the intent, a set of formal properties or attributes, a set of frames.

E_2 the members of a VerbNet class, a set of English verbs and finally

I_2 a set of frames describing the class.

asd is a similarity function defined on the set of frames and asd_{vn} a similarity function on French and English verbs (eg. as defined in Definition 17).

$$\begin{aligned}
Sim((E_1, I_1), (E_2, I_2)) := & \left[\frac{1}{\max(|E_1|, |E_2|)} \max_{P \in \mathcal{P}(E_1, E_2)} \left(\sum_{(a,b) \in P} asd_{vn}(a, b) \right) \right] * w \\
& + \left[\frac{1}{\max(|I_1|, |I_2|)} \max_{P \in \mathcal{P}(I_1, I_2)} \left(\sum_{(a,b) \in P} asd(a, b) \right) \right] * (1 - w)
\end{aligned} \tag{3}$$

where

\mathcal{P} is the set of matchings of the sets E_1, E_2 and I_1, I_2 respectively.
 $0 \leq w \leq 1$ is a weight.

Since in our example the frames can be directly matched, we can simplify the definition to:

Definition 19 (Similarity of a formal concept and a VerbNet class, simplified) Let

E_1 be the extent of the formal concept, a set of French verbs,
 I_1 the intent, a set of formal properties or attributes, a set of frames,
 E_2 members of the VerbNet class, a set of English verbs and finally
 I_2 a set of frames describing the class.

Let asd_{vn} be a similarity function on French and English verbs (eg. as defined in Definition 17).

$$\begin{aligned}
Sim((E_1, I_1), (E_2, I_2)) := & \left[\frac{1}{\max(|E_1|, |E_2|)} \max_{P \in \mathcal{P}(E_1, E_2)} \left(\sum_{(a,b) \in P} asd_{vn}(a, b) \right) \right] * w \\
& + \frac{|I_1 \cap I_2|}{\max(|I_1|, |I_2|)} * (1 - w)
\end{aligned} \tag{4}$$

where

\mathcal{P} is the set of matchings of the sets E_1, E_2 .
 $0 \leq w \leq 1$ is a weight.

We are now ready to test this method on the running example described in the following section.

7.2 The Running Example

We will use the following set of French verbs: *briser, fendre, casser, couper gratter, érafler, frapper, battre, cogner, toucher, piquer* and *embrasser*, and the following set of VerbNet frames and alternations:

| | |
|------------------|---------------------------|
| BASIC TRANSITIVE | ↔ INSTRUMENT SUBJECT |
| BASIC TRANSITIVE | ↔ MIDDLE |
| BASIC TRANSITIVE | ↔ TRANS BODY-PART OR REFL |
| BASIC TRANSITIVE | ↔ INTRANSITIVE INCHOATIVE |
| BASIC TRANSITIVE | ↔ CONATIVE |

| | INSTRUMENT SUBJECT | MIDDLE | TRANS. BODY-PART OR REFL. | INTRANSITIVE INCHOATIVE | CONATIVE |
|-----------|--------------------|--------|---------------------------|-------------------------|----------|
| briser | x | x | | x | |
| fendre | x | x | | x | |
| casser | x | x | | x | |
| couper | x | x | x | | x |
| gratter | x | x | x | | x |
| érafler | x | x | x | | x |
| frapper | | | x | | x |
| battre | | | x | | x |
| cogner | | | x | | x |
| toucher | | | x | | |
| piquer | | | x | | |
| embrasser | | | x | | |

Table 10: Formal context for 12 French verbs from the break, cut, hit and touch class.

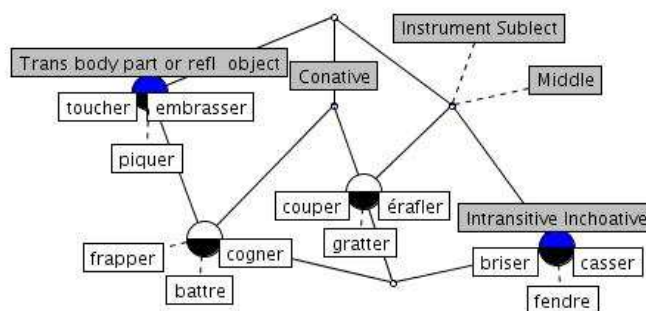


Figure 6: The concept lattice for the formal context of Table 10.

As each of the above alternations reflect a variation of the BASIC TRANSITIVE syntactic frame, we will use only the second alternation component (set in bold face above) as attribute identifier.

By specifying the frames accepted by each of these verbs we obtain the formal context depicted in Table 10.

Applying formal concept analysis we obtain the concept lattice in Figure 6.

As a result, we obtain the verb classes (concepts) given in Table 11.

We want to investigate how these classes relate to the VerbNet classes *break-45.1*, *touch-20.1*, *hit-18.1-1* and *cut-21.1-1*. When computing the attribute similarity scores, we take into account all the VerbNet frames of these classes, regardless of their occurrence in the example.

| | | |
|------------|-----------------------------------|---|
| c.1 | <i>piquer, toucher, embrasser</i> | TRANS BODY PART OR REFLEXIVE OBJECT |
| c.4 | <i>cogner, frapper, battre</i> | TRANS BODY PART OR REFLEXIVE OBJECT, CONATIVE |
| c.5 | <i>briser, fendre, casser</i> | INSTRUMENT SUBJECT, MIDDLE, INTRANSITIVE INCHOATIVE |
| c.6 | <i>gratter, erafler, couper</i> | CONATIVE, INSTRUMENT SUBJECT, MIDDLE |

Table 11: Computed verb classes for the running example.

7.3 Concept Similarity for the Running Example

Thus, to match a formal concept to a VerbNet class we have to compute similarity scores for each pair of a concept and a class and choose the pair with maximal score.

Table 12 shows an excerpt of the similarity scores between the English and French verbs of our example (see Table 19 in the Appendix for the entire table).

| | <i>battre</i> | <i>briser</i> | <i>casser</i> | <i>cogner</i> | <i>couper</i> | <i>embrasser</i> | <i>erafler</i> | <i>fendre</i> | <i>frapper</i> | <i>gratter</i> | <i>piquer</i> | <i>toucher</i> |
|---------|---------------|---------------|---------------|---------------|---------------|------------------|----------------|---------------|----------------|----------------|---------------|----------------|
| bang | 0.17 | 0 | 0.17 | 0 | 0 | 0 | 0 | 0 | 0.17 | 0 | 0 | 0 |
| bash | 0 | 0 | 0 | 0.20 | 0 | 0 | 0 | 0 | 0.20 | 0 | 0 | 0 |
| break | 0 | 0.17 | 0.17 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| cut | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| graze | 0 | 0 | 0 | 0 | 0 | 0 | 0.14 | 0 | 0 | 0 | 0 | 0 |
| hew | 0 | 0 | 0 | 0 | 0.25 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| hit | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| scratch | 0 | 0 | 0 | 0 | 0 | 0 | 0.09 | 0 | 0 | 0.09 | 0 | 0 |
| shatter | 0 | 0.17 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| strike | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.33 | 0 | 0 | 0 |
| touch | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.25 |
| whack | 0.50 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

Table 12: Similarity scores for some French and English verbs for the running example.

The scores were computed automatically as follows: We translated the verb members of the involved VerbNet classes by using dictionaries available on-line²². If one of the French verbs occurred in the set of found translations we assigned the pair a score of 1. We obtained the final score displayed in Table 12 by dividing the initial score by the number of found translations. For instance, the score of 0.17 for *bang* and *battre* is obtained as follows: looking up *bang* in the dictionaries we obtain 6 French verbs, one of which is *battre*. Therefore the similarity score assigned to the pair *bang* and *battre* will be of 1/6.

Next we compute the similarity scores between a concept and a VerbNet class. Concepts and VerbNet classes each consist of a set of member verbs and a set of frames. The similarity between a concept and a VerbNet class is computed by combining the similarity scores between its sets of verbs and its sets of attributes following Definition 19.

We obtain the attribute similarity score for concept c_G and VerbNet class c_{VN} by computing $\frac{|I_1 \cap I_2|}{\max(|I_1|, |I_2|)}$ where I_1 is the frame set of c_G and I_2 the frame set of c_{VN} .

We compute the object similarity for c_G and c_{VN} by solving the following maximum weight matching problem: we consider the graph whose vertexes are elements of the set E_1 , the set of verbs of c_G , and of the set E_2 , the set of verbs of c_{VN} . A French verb v_{fr} of E_1 and an English verb v_{en} of E_2 are connected if $Sim(v_{fr}, v_{en}) > 0$ (from

²²<http://www.wordreference.com/fren/> and <http://dictionnaire.mediadico.com/traduction/dictionnaire.asp/anglais-francais/>

Table 12). The connecting edge is weighted by the similarity score $Sim(v_{fr}, v_{en})$ from the table.

Figure 7 illustrates the above setting for a set E_2 of English verbs: *bash*, *bang*, *strike*, *whack*, *break*, *shatter* and a set E_1 of French verbs: *cogner*, *frapper*, *battre*, *casser*, *briser*, *couper*. Figure 7(a) shows the full graph obtained from the verb sets and the similarity Table 12 and Figure 7(b) displays an arbitrary matching.

Next we consider all matchings of this graph (cf. Definition 15) and select one where the sum of the edge weights is maximal. This matching is a solution to this maximum weight matching problem and the maximal sum of weights is the object similarity score between c_G and c_{VN} ²³.

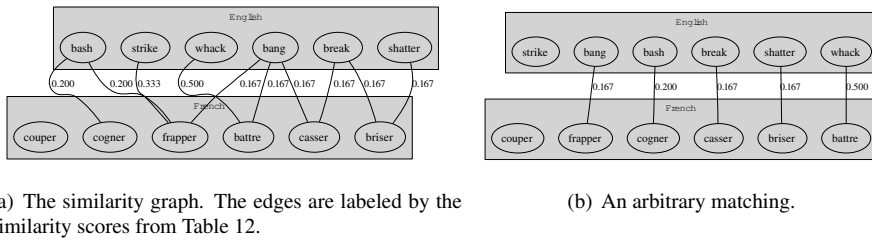


Figure 7: The maximum weight matching problem for a set of French and English verbs of the running example. Fig. (a) shows the full graph and Fig. (b) an arbitrary matching. A solution is a matching where the sum of the edge weights is maximal.

Finally, the attribute, object and overall similarity scores for the computed concepts and considered VerbNet classes are depicted in Table 13²⁴.

| | break-45.1 | | | cut-21.1-1 | | | hit-18.1-1 | | | touch-20-1 | | |
|-----|-------------|-------------|-------------|-------------|------|-------------|------------|-------------|---------|-------------|-------------|-------------|
| | att | obj | overall | att | obj | overall | att | obj | overall | att | obj | overall |
| c_1 | 0.00 | 0.00 | 0 | 0.08 | 0.00 | 0.04 | 0.09 | 0.00 | 0.05 | 0.33 | 0.17 | 0.25 |
| c_4 | 0.00 | 0.00 | 0 | 0.17 | 0.00 | 0.08 | 0.18 | 0.14 | 0.16 | 0.33 | 0.11 | 0.22 |
| c_5 | 0.33 | 0.04 | 0.18 | 0.17 | 0.00 | 0.08 | 0.09 | 0.01 | 0.05 | 0.00 | 0.00 | 0 |
| c_6 | 0.22 | 0.00 | 0.11 | 0.25 | 0.03 | 0.14 | 0.18 | 0.00 | 0.09 | 0.00 | 0.05 | 0.02 |

Table 13: Similarity Scores for Computed and VerbNet classes for the running example. The weight is 0.5. The classes are as follows: **c_1**: *piquer*, *toucher*, *embrasser*, **c_4**: *cogner*, *frapper*, *battre*, **c_5**: *briser*, *fendre*, *casser*, **c_6**: *gratter*, *erafler*, *couper*. Maximum values are set in bold face.

By the overall similarity score the system correctly assigns the classes c_1 , c_5 and c_6 to the VerbNet classes *touch-20-1*, *break-45.1* and *cut-21.1-1*. The class c_4 is wrongly assigned to the *touch-20-1* VerbNet class, whereas the object similarity score correctly related this class to *hit-18.1-1*. However, in the case of c_6 the object similarity score selects the *touch-20-1* class instead of the expected *cut-21.1-1* but the overall similarity is balanced by the attribute similarity score. When adjusting the weight to 0.9, i.e. reinforcing the weight of the object similarity, all computed classes are related to the expected VerbNet classes (cf. Table 14).

²³In fact the weight sum must be divided by the maximum cardinality of E_1 and E_2

²⁴The scores were computed using `Algorithm::Munkres`, a Perl implementation of the Hungarian Algorithm for solving the so-called *assignment problem* in polynomial time (cf. [Mun57] and [BL71]). The assignment problem is equivalent to the maximum weight matching problem

| | break-45.1 | cut-21.1-1 | hit-18.1-1 | touch-20-1 |
|-----|-------------------|-------------------|-------------------|-------------------|
| c.1 | 0 | 0.01 | 0.01 | 0.18 |
| c.4 | 0 | 0.02 | 0.14 | 0.13 |
| c.5 | 0.07 | 0.02 | 0.02 | 0 |
| c.6 | 0.02 | 0.05 | 0.02 | 0.04 |

Table 14: Overall similarity scores for computed and VerbNet classes for the running example. The weight is 0.9. Classes are as above (Table 13). Maximum values are set in bold face.

Although these results may seem promising, they need to be used with caution for different reasons: first of all we do not know to what extent (if at all) French verb classes relate to English verb classes. Second, the way we establish similarity scores between French and English verbs depends strongly on the quality of the translation and is therefore coarse and noisy. The translation quality can only be poor because of the notoriously high ambiguity. Disambiguating the translations, for example by exploiting available WordNet sense mappings²⁵ or the translations of Dicovalence (cf. 4.3) would certainly improve the scores. Last but not least the example is too small and thus only partly conclusive.

On the other hand the example shows that these techniques are applicable and operational.

7.4 Comparing Class Hierarchies

In Section 6.3.1 we saw that classifications need to be compared both class-wise and as hierarchical structures. We also found (Section 6.3.1) that a comparison of the hierarchical structure did not (yet) make sense in our use case. Nevertheless, the technique described in the previous section (Sec. 7.1) gives rise to a method for comparing class hierarchies. In the following we outline the basic idea.

In Section 7.1 we compute a similarity score for each pair of a concept and a VerbNet class. This translates to the “Graph maximum weighted matching” problem as follows: We have to find the matching between concepts and VerbNet classes, where the sum of similarity scores is maximal. This maximal matching gives a concept \leftrightarrow VerbNet class matching optimised on the classification level (rather than the single class level) and the maximum score measures the classification similarity. As we can easily convince ourselves by means of Fig. 7, this is in fact very similar to Definition 16: We solve the “maximum weighted matching” problem on pairs of concepts and classes instead of the pairs of verbs in Definition 16.

8 Conclusion

In this work we used *Formal Concept Analysis* to automatically acquire syntactic-semantic verb class and diathesis alternation information from three French syntax lexicons.

We created large scale verb classes and compared their verb and frame distribution. The acquired verb classifications are unsatisfactory for reasons related first to the quantity and quality of the data and second to the naïve FCA application strategy. Nonetheless, they give many clues to improvements and directions for further investigations and help augment the knowledge about the task of semantic verb classification in general and the resources in particular.

²⁵VerbNet maps many verbs to their WordNet senses. For French, ([SF08]) automatically built a French WordNet-like resource based on multilingual parallel corpora.

We extracted pairs of argument frames from the French resources, exposing frequency and co-occurrence information, thus providing more sound grounds for linguist experts to judge their adequacy as French diathesis alternations.

Finally, we addressed the evaluation issue by proposing an evaluation methodology of the acquired classifications with regard to the English VerbNet and implementing it on a small example.

Directions for further work are manifold. They centre around two themes:

1. Improving and extending the resources and the FCA application technique.
2. Evaluation: developing the method discussed in Sec. 7.1, exploring further strategies.

Theme 1 comprises the conclusion of the commenced validation and normalisation task, which will allow for an extension of the coverage by fusing the available resources. This work goes hand in hand with a refinement of the FCA application techniques.

Wrt. Theme 2, the next obvious enterprise is to test whether the proposed method scales. Furthermore, we saw that this method strongly relies on a similarity measure between French and English verbs. Much work has been focused on translation in contexts or aligned corpora and more in depth investigations in this area could lead to better similarity measures.

On the other hand, there exists a variety of evaluation measures for classifications from diverse areas such as web-page clustering and ontology learning and mapping which could also be explored. These would need to be tested on a “representative” sample of the classification, which is not available currently. Therefore, strategies to identify “representative” parts of a concept lattice and of the data would also be beneficial.

Last but not least, we saw that an important feature of our targeted classification is the semantic coherence of its verb classes. As currently we know of no automatic means of its assessment, this is also an interesting direction for further investigations.

Ultimately we believe that a comprehensive evaluation will not be possible without a Gold Standard. As mentioned previously there is a long way to go towards its creation but various techniques from the literature, some using FCA or ontology learning, show ways towards the acquisition of the required reference.

Acknowledgements

First I would like to thank my supervisors, Claire Gardent and Fabienne Venant, who were always available for advice and whose guidance was always demanding and encouraging but never constricting. I am also grateful to Yannick Toussaint, Amine Mohamed Rouane Hacene, Philipp Cimiano and Jens Gustedt for discussions which provided me with many helpful insights.

Appendix

A Tables

| | | |
|------------|---------|---|
| Volem | Members | rôder revenir régner rebondir passer partir flotter camper affluer |
| | Intent | 100 15 191 anti_np caus_np_pp |
| Ladl | Members | sourire souffrir jouir intervenir dîner dépérir déjeuner |
| | Intent | 100 16 180 anti_np caus_np_pp |
| Dicovalece | Members | évaluer savourer redouter photographier mésestimer détecter boycotter |
| | Intent | 32R2 6 |
| VerbNet | Members | vendanger râteler poncer limer exciser dépiauter défricher décaper cureter |
| | Intent | 37E 38LS |
| Dicovalece | Members | téléphoner télégraphier réciter réaffirmer raconter prescrire narrer gaspiller |
| | Intent | P0 (P1) (P2) |
| VerbNet | Members | voiler repeupler libérer fleurir épurer embellir dissocier couronner blanchir alourdir accuser |
| | Intent | P0, P0 P1 (P3), P0 P1 |
| VerbNet | Members | broadcast, cable, e-mail, fax, modem, netmail, phone, radio, relay, satellite, semaphore, sign, signal, telecast, telegraph, telephone, telex, wire, wireless |
| | Frames | Basic Transitive, Dative, FOR-TO-INF, NP-PP Topic-PP, NP-PP to-PP, NP-QUOT Recipient Object, NP-S Recipient Object, NP-TO-INF-AC Recipient Object, NP-WHEN-TO-INF Recipient Object, . . . |
| VerbNet | Members | aromatize, asphalt, bait, black, board, bread, brick, bridle, bronze, butter, buttonhole, cap, . . . (115 members) |
| | Frames | NP-PP Theme-PP, Transitive Destination Object |

Table 15: Examples of verb classes for Volem, Ladl and Dicovalece The first row for each entry shows the verbs, the second the corresponding intent. The last entry shows sample classes and frames of VerbNet

| | caus_refl_pr_np_pp | pas_etre_part_np_pp | 102 | state_2np | caus_2np | anti_pr_np | 162 | 141 | 16 | 100 | 190 |
|------------|--------------------|---------------------|-----|-----------|----------|------------|-----|-----|----|-----|-----|
| abaïsser | | x | | | x | x | | | | | |
| abandonner | | x | x | | x | x | x | | | | x |
| abattre | | x | | | x | x | | | | | |
| abonner | x | x | | | | | | | | | |
| aborder | | | | | x | x | x | x | | | |
| aboutir | | | | | | | | | x | | |
| abriter | x | x | x | | x | x | x | | | | |
| abroger | | x | | | x | x | x | | | | |
| abré ger | | x | | | x | x | | | | x | |
| abî mer | x | x | x | x | x | x | | x | | | |

Table 16: Volem. The first 10 rows of the Volem context (and 10 attributes)

| | battre | briser | casser | cogner | couper | embrasser | fler | fendre | frapper | gratter | piquer | toucher |
|--------|--------|--------|--------|--------|--------|-----------|------|--------|---------|---------|--------|---------|
| bang | 0.17 | 0 | 0.17 | 0 | 0 | 0 | 0 | 0 | 0.17 | 0 | 0 | 0 |
| bash | 0 | 0 | 0 | 0.20 | 0 | 0 | 0 | 0 | 0.20 | 0 | 0 | 0 |
| batter | 1.00 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| beat | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| break | 0 | 0.17 | 0.17 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| bump | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| butt | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| caress | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

Continued on next page

| | 37E | 7 | 32C | 32H | 11 | 37M6 | 32R3 | 38L | 6 | 31H | 32RA |
|-------------|-----|---|-----|-----|----|------|------|-----|---|-----|------|
| abaïsser | | | | | x | | | x | | | x |
| abandonner | | x | | x | | | | | x | x | |
| abasourdir | | | | x | | | | | | | x |
| abattre | x | | | x | | | x | | | | |
| abdiquer | | | | | | | x | | | | |
| abhorrer | | | | | | | | | | | |
| abjurer | | | | | | | x | | | | |
| ablutionner | | | | | | x | | | | | |
| abê tir | | | | | x | | | | | | |
| abî mer | | | x | | | | | | | | |

Table 17: **Ladl**. The first 10 rows of the Ladl context (and 10 attributes)

| | P0 (PP < en faveur de >) | P0 (PL) | P0 PL | P0 (P2) | P0 (PP < devant >) | P0 (P1) | EMPTY | P0 PP < en > | P0 | P0 (PP < vers >) (PQ) | P0 P2 |
|------------|---------------------------|---------|-------|---------|---------------------|---------|-------|--------------|----|------------------------|-------|
| abaïsser | | | x | | x | x | | | x | x | x |
| abandonner | | x | | x | | x | | | | | |
| abattre | | | | | | | | | | | |
| abdiquer | x | | | | | | | | | | |
| abhorrer | | | | | | | x | | | | |
| abjurer | | | | | | x | | | | | |
| abolir | | | | | | | | | | | |
| abominer | | | | | | | | | | | |
| abonder | | x | | | | | | x | | | |
| abî mer | | | | | | | | | x | | |

Table 18: **Dicovalence**. The first 10 rows of the Dicovalence context (and 10 attributes)

| | battre | briser | casser | cogner | couper | embrasser | fler | fendre | frapper | gratter | piquer | toucher |
|----------|--------|--------|--------|--------|--------|-----------|------|--------|---------|---------|--------|---------|
| chip | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| click | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| clip | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| crack | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| crash | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| crush | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| cut | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| dash | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| drum | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| fondle | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| fracture | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| graze | 0 | 0 | 0 | 0 | 0 | 0 | 0.14 | 0 | 0 | 0 | 0 | 0 |
| hack | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| hammer | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| hew | 0 | 0 | 0 | 0 | 0.25 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| hit | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| kick | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| kiss | 0 | 0 | 0 | 0 | 0 | 0.25 | 0 | 0 | 0 | 0 | 0 | 0 |

Continued on next page

| | battre | briser | casser | cogner | couper | embrasser | fler | fendre | frapper | gratter | piquer | toucher |
|----------|--------|--------|--------|--------|--------|-----------|------|--------|---------|---------|--------|---------|
| knock | 0 | 0 | 0 | 0.10 | 0 | 0 | 0 | 0 | 0.10 | 0 | 0 | 0 |
| lash | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| lick | 0.33 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| nudge | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| pat | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| peck | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| pinch | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| pound | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| prod | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| rap | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.14 | 0 | 0 | 0 |
| rip | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| saw | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| scrape | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| scratch | 0 | 0 | 0 | 0 | 0 | 0 | 0.09 | 0 | 0 | 0.09 | 0 | 0 |
| shatter | 0 | 0.17 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| slap | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| slash | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| smack | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| smash | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| snap | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| snip | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| splinter | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| split | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| sting | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| strike | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.33 | 0 | 0 | 0 |
| stroke | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| tamp | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| tap | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| tear | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| thump | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| thwack | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| tickle | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| touch | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.25 |
| whack | 0.50 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

Table 19: English – French similarity scores. Similarity scores between English and French verbs for the running example. The rows contain the scores for the verbs of the considered VerbNet classes, the columns the considered French verbs.

B Implementation

For this work we used

1. tools to compute FCA related devices,
2. tools to compute solutions for the maximum weight matching problem and
3. experimentation tools.

FCA tools We used the following FCA software to compute concepts, lattices, association rules, their support and confidence:

Galicia <http://www.iro.umontreal.ca/~galicia/>, implemented in Java

ConExp We used *The Concept Explorer*²⁶, mainly as a visualisation tool.

Colibri-concepts is a command line tool for formal concept analysis, implemented in C²⁷.

FCA::Context A self-made Perl implementation to be used on small examples, which we plan to release on CPAN²⁸ as soon as it is sufficiently tested and documented.

The Maximum Weight Matching problem We used `Algorithm::Munkres`²⁹, a Perl implementation of the Munkres (aka Hungarian) algorithm which computes solutions to the *assignment problem* in polynomial time (cf. [Mun57] and [BL71]). The assignment problem is equivalent to the maximum weight matching problem. The implementation was sufficient for our small example, but it is not clear whether and how it would scale up. However, there exist several other software packages which implement this algorithm in C.

Experimentation tools All experiments were conducted using Perl scripts. For visualisation we also used the Aduna software³⁰.

We registered a project on the InriaGforge³¹ and used its collaborative development services. All developed resources, in terms of software, files and documentation are available there on demand.

²⁶<http://conexp.sourceforge.net/>, implemented in Java

²⁷<http://code.google.com/p/colibri-concepts/>

²⁸The Comprehensive Perl Archive Network, <http://www.cpan.org>.

²⁹<http://search.cpan.org/~tpederse/Algorithm-Munkres-0.07/>

³⁰cf. [FSvH05], <http://www.aduna-software.com/home/overview.view>

³¹<http://gforge.inria.fr/projects/verbnet-fr/>

References

- [BGL76] J.-P. Boons, A. Guillet, and C. Leclère. *La structure des phrases simples en français. I : Constructions intransitives*. Droz, Genève, 1976.
- [BL71] François Bourgeois and Jean-Claude Lassalle. An Extension of the Munkres Algorithm for the Assignment Problem to Rectangular Matrices. *Commun. ACM*, 14(12):802–804, 1971.
- [BM70] M. Barbut and B. Monjardet, editors. *L'Ordre et la classification*. Algèbre et combinatoire, tome II. Hachette, 1970.
- [BTP⁺02] Y. Bastide, R. Taouil, N. Pasquier, G. Stumme, and L. Lakhal. Pascal: un algorithme d'extraction des motifs frents. *Technique et Science Informatiques (TSI)*, 21(1):65–95, 2002.
- [CHS05] Philipp Cimiano, Andreas Hotho, and Steffen Staab. Learning concept hierarchies from text corpora using formal concept analysis. *Journal of Artificial Intelligence Research (JAIR)*, 24:305–339, AUG 2005.
- [Del05] Klaas Dellschaft. Measuring the similarity of concept hierarchies and its influence on the evaluation of learning procedures. Master's thesis, Universität Koblenz Landau, Campus Koblenz, Fachbereich 4 Informatik, Institut für Computervisualisirk, 2005.
- [DS06] Klaas Dellschaft and Steffen Staab. On how to perform a gold standard based evaluation of ontology learning. In I. Cruz et al., editor, *Proceedings of the 5th International Semantic Web Conference (ISWC)*, LNCS 4273, pages 228–241. Springer Verlag, 2006.
- [Fel98] C. Fellbaum. *WordNet: An Electronical Lexical Database*. The MIT Press, Cambridge, MA, 1998.
- [For06] Anna Formica. Ontology-based concept similarity in formal concept analysis. *Information Sciences*, 176(18):2624–2641, September 2006.
- [For08] Anna Formica. Concept Similarity in Formal Concept Analysis: An Information Content Approach. *Know.-Based Syst.*, 21(1):80–87, 2008.
- [FSvH05] Christiaan Fluit, Marta Sabou, and Frank van Harmelen. Ontology-based information visualisation: Towards semantic web applications. In Vladimir Geroimenko, editor, *Visualising the Semantic Web (2nd edition)*. Springer Verlag, 2005.
- [FVSD⁺02] Ana Fernandez, Gloria Vazquez, Patrick Saint-Dizier, Farah Benamara, and Mouna Kamel. The volem project: a framework for the construction of advanced multilingual lexicons. In *LEC '02: Proceedings of the Language Engineering Conference (LEC'02)*, page 89, Washington, DC, USA, 2002. IEEE Computer Society.
- [Gal86] Zvi Galil. Efficient algorithms for finding maximum matching in graphs. *ACM Comput. Surv.*, 18(1):23–38, 1986.
- [GGPF05] C. Gardent, B. Guillaume, G. Perrier, and I. Falk. Extracting subcategorisation information from Maurice Gross' Grammar Lexicon. *Archives of Control Sciences*, 15(LI):253–264, 2005.

- [GGPF06] Claire Gardent, Bruno Guillaume, Guy Perrier, and Ingrid Falk. Extraction d'information de sous-catégorisation à partir des tables du LADL. In *Traitement Automatique de la Langue Naturelle - TALN 2006 Actes de la 13ème conférence sur le Traitement Automatique de la Langue Naturelle*, Leuven/Belgique, 04 2006.
- [GL92] A. Guillet and C. Leclère. *La structure des phrases simples en français. Constructions transitives locatives*. Droz, Genève, 1992.
- [Gro75] M. Gross. *Méthodes en syntaxe*. Hermann, 1975.
- [GS03] Bernhard Ganter and Gerd Stumme. Formal Concept Analysis: Methods and Applications in Computer Science, 2003. Lecture Notes, Otto-von-Guericke-Universität Magdeburg.
- [HNHRV07] Marianne Huchard, Amedeo Napoli, M. Hacene Rouane, and Petko Valtchev. Mining Description Logics Concepts With Relational Concept Analysis. In Patrice Bertrand Paula Brito, Guy Cucumel and Francisco de Carvalho, editors, *Selected Contributions in Data Analysis and Classification*, Studies in Classification, Data Analysis, and Knowledge Organization, pages 259–270. Springer Berlin Heidelberg, 08 2007.
- [Joa02] Eric Joanis. Automatic verb classification using a general feature space. Master's thesis, Department of Computer Science, University of Toronto, October 2002.
- [JS03] Eric Joanis and Suzanne Stevenson. A general feature space for automatic verb classification. In *Proceedings of the 11th Conference of the European Chapter of the Association for Computational Linguistics (EACL-03)*, Budapest, Hungary, April 2003.
- [JSJ06] Eric Joanis, Suzanne Stevenson, and D. James. A general feature space for automatic verb classification. *Natural Language Engineering*, 2006.
- [KB04] Anna Korhonen and Ted Briscoe. Extended lexical-semantic classification of english verbs. In Dan Moldovan and Roxana Girju, editors, *HLT-NAACL 2004: Workshop on Computational Lexical Semantics*, pages 38–45, Boston, Massachusetts, USA, May 2 - May 7 2004. Association for Computational Linguistics.
- [KKRP06] Karin Kipper, Anna Korhonen, Neville Ryant, and Martha Palmer. Extending verbnet with novel verb classes. In *Proceedings of 5th international conference on Language Resources and Evaluation*, 2006.
- [Lev93] Beth Levin. *English Verb Classes and Alternations: a preliminary investigation*. University of Chicago Press, Chicago and London, 1993.
- [Lin98] Dekang Lin. An Information-Theoretic Definition of Similarity. In *ICML '98: Proceedings of the Fifteenth International Conference on Machine Learning*, pages 296–304, San Francisco, CA, USA, 1998. Morgan Kaufmann Publishers Inc.
- [LK08] Anke Lüdeling and Merja Kytö, editors. *Corpus Linguistics. An International Handbook*. Handbooks of Linguistics and Communication Science. Mouton de Gruyter, Berlin, 2008. To appear.

- [Mao07] Ming Mao. Ontology mapping: An information retrieval and interactive activation network based approach. In *ISWC/ASWC*, pages 931–935, 2007.
- [MPS07] Ming Mao, Yefei Peng, and Michael Spring. A profile propagation and information retrieval based ontology mapping approach. In *SKG '07: Proceedings of the Third International Conference on Semantics, Knowledge and Grid*, pages 164–169, Washington, DC, USA, 2007. IEEE Computer Society.
- [MS01] Paola Merlo and Suzanne Stevenson. Automatic verb classification based on statistical distributions of argument structure. *Computational Linguistics*, 27(3):373–408, 2001.
- [MS02] Alexander Maedche and Steffen Staab. Measuring similarity between ontologies. In *EKAW '02: Proceedings of the 13th International Conference on Knowledge Engineering and Knowledge Management. Ontologies and the Semantic Web*, pages 251–263, London, UK, 2002. Springer-Verlag.
- [Mun57] J. Munkres. Algorithms for the Assignment and Transportation Problems. *Journal of the Society of Industrial and Applied Mathematics*, 5(1):32–38, March 1957.
- [Res95] Philip Resnik. Using Information Content to Evaluate Semantic Similarity in a Taxonomy. In *IJCAI*, pages 448–453, 1995.
- [Sch03] Sabine Schulte im Walde. *Experiments on the Automatic Induction of German Semantic Verb Classes*. PhD thesis, Institut für Maschinelle Sprachverarbeitung, Universität Stuttgart, 2003. Published as AIMS Report 9(2).
- [Sch05] Karin Kipper Schuler. *Verbnet: A Broad-Coverage, Comprehensive Verb Lexicon*. PhD thesis, Faculties of Computer and Information Science of the University of Pennsylvania, 2005.
- [Sch06] Sabine Schulte im Walde. Experiments on the automatic induction of german semantic verb classes. *Computational Linguistics*, 32(2):159–194, 2006.
- [Sch08] Sabine Schulte im Walde. The Induction of Verb Frames and Verb Classes from Corpora. In Lüdeling and Kytö [LK08], chapter 61. To appear.
- [SD95] Patrick Saint-Dizier. Verb semantic classes in french. version 1. Technical report, IRIT – CNRS, Toulouse, 1995.
- [SD96] Patrick Saint-Dizier. Constructing verb semantic classes for french: Methods and evaluation. In *COLING*, pages 1127–1130, 1996.
- [SD99] P. Saint-Dizier. Alternation and verb semantic classes for french: Analysis and class formation. In *Predicative forms in natural language and in lexical knowledge bases*. Kluwer Academic Publishers, 1999.

- [SF08] Benoît Sagot and Darja Fišer. Building a Free French Wordnet from Multilingual Resources. In *OntoLex 2008*, Marrakech, Morocco, 2008.
- [Spo02] Caroline Sporleder. A galois lattice based approach to lexical inheritance hierarchy learning. In *OLT '02: Proceedings of the ECAI Workshop on Machine Learning and Natural Language Processing for Ontology Engineering*, Lyon, France, 2002.
- [STB⁺02] G. Stumme, R. Taouil, Y. Bastide, N. Pasquier, and L. Lakhal. Computing iceberg concept lattices with Titanic. *J. Knowledge and Data Engineering (KDE)*, 42(2):189–222, 2002.
- [SWGM05] Marta Sabou, Chris Wroe, Carole Goble, and Gilad Mishne. Learning domain ontologies for web service descriptions: an experiment in bioinformatics. In *WWW '05: Proceedings of the 14th international conference on World Wide Web*, pages 190–198, New York, NY, USA, 2005. ACM.
- [VB98] P. Vossen and L. Bloksma. Categories and Classifications in EuroWordNet. In R. Catro A. Rubio, N. Gallardo and A. Tejada, editors, *Proceedings of First International Conference on Language Resources and Evaluation*, Granada, May 1998.
- [Vos98] Piek Vossen. EuroWordNet: Building a Multilingual Database with Wordnets for European Languages. In M. Nilsson K. Choukri, D. Fry, editor, *The ELRA Newsletter*, volume 3. ELRA, 1998.
- [Wil82] Rudolf Wille. Restructuring Lattice Theory: an Approach Based on Hierarchies of Concepts. In Ivan Rival, editor, *Ordered sets*, pages 445–470, Dordrecht–Boston, 1982. Reidel.