

PAROLE/SIMPLE ‘Lemon’ ontology and lexicons

Marta Villegas and Núria Bel
Universitat Pompeu Fabra

Abstract. The PAROLE/SIMPLE ‘Lemon’ Ontology and Lexicon are the OWL/RDF version of the PAROLE/SIMPLE lexicons (defined during the PAROLE (LE2-4017) and SIMPLE (LE4-8346) IV FP EU projects) once mapped onto Lemon model and LexInfo ontology. Original PAROLE/SIMPLE lexicons contain morphological, syntactic and semantic information, organized according to a common model and to common linguistic specifications for 12 European languages. The data set we describe includes the PAROLE/SIMPLE model mapped to Lemon and LexInfo ontology and the Spanish & Catalan lexicons. All data are published in the Data Hub and are distributed under CC Attribution 3.0 Unported license. The Spanish lexicon contains 199466 triples and 7572 lexical entries fully annotated with syntactic and semantic information. The Catalan lexicon contains 343714 triples and 20545 lexical entries annotated with syntactic information half of which are also annotated with semantic information. In this paper we describe the resulting data, the mapping process and the benefits obtained. We demonstrate that the Linked Open Data principles prove essential for datasets such as original PAROLE/SIMPLE lexicons where harmonization and interoperability was crucial. The resulting data is lighter and better suited for exploitation. In addition, it eases further extensions and links to external resources such as WordNet, lemonUby, DBpedia etc.

Keywords: lexicon, ontology, open linked data, RDF, OWL, LE-PAROLE, SIMPLE, LexInfo, Lemon

1. Introduction

The PAROLE/SIMPLE ‘Lemon’ Ontology is the OWL/RDF version of the PAROLE & SIMPLE lexicon models (defined during the PAROLE LE2-4017 and SIMPLE LE4-8346 projects) once mapped to Lemon¹ and LexInfo² models.

1.1. PAROLE/SIMPLE lexicons

Original PAROLE/SIMPLE lexicons contain morphological, syntactic and semantic information organized according to a common model and to common linguistic specifications. PAROLE was the first project producing corpora and lexicons in so many languages³ and built according to the same design principles, linguistic specifications and representation format. The model was based on EAGLES recommendations for morphosyntactic information and verb syntax [7] and on the extended GENELEX model [1].

¹ <http://lemon-model.net/>

² <http://lexinfo.net/>

³ Catalan, Danish, Dutch, English, Finnish, French, German, Greek, Italian, Portuguese, Spanish and Swedish

The goal of SIMPLE project was to add semantic information to the set of harmonized multifunctional lexicons built for 12 European languages by the PAROLE consortium. All PAROLE/SIMPLE lexicons were defined against a common model defined in the DTD. Thus all PAROLE/SIMPLE lexicons are XML files valid against the same DTD⁴. In addition, a good number of ‘descriptive’ elements were defined and shared by all SIMPLE lexicons. Essentially, these include: (i) Template assignment: meant to guarantee coherent encoding, across sites and languages, (ii) Domain information, (iii) Semantic class information, (iv) Semantic features: distinctive features used to better specify the semantic class of a sense, and for the definition of selectional restrictions on the arguments (v) Semantic Roles and (vi) Semantic Relations.

1.2. LMF, Lemon and LexInfo

LMF [5] (Lexical Markup Framework) is an ISO standard (ISO-24613:2008) for Natural Language

⁴ Original PAROLE/SIMPLE lexicons were in SGML so we previously converted them into XML.

Processing lexicons. LMF combines the best designs and methods from many existing NLP lexicons⁵. LMF models are represented by UML classes, associations among the classes, and a set of ISO-12620 data categories that function as UML attribute-value pairs. LMF includes an XML DTD where XML elements in the DTD are transcoded from the UML class diagrams. The class adornment is implemented as a set of *feat* elements

Lemon [6] ('lexicon model for ontologies' developed by the Monnet project <http://www.monnet-project.eu/>) is a model for modeling lexicon based on LMF and expressed in RDF. The Lemon model consists of a core path defined as: *OntologyEntity* ↔ *LexicalSense* ↔ *LexicalEntry* → *LexicalForm* → *WrittenRepresentation*. Lemon is highly compliant with LMF.

LexInfo [1] is a model for the linguistic grounding of ontologies and as such allows for the association linguistic information (such as part-of-speech, subcategorization frames etc.) with ontology elements (such as concepts, relations, individuals, etc.). LexInfo builds on the Lemon model and it is also highly compliant with LMF.

1.3. The mapping

Mapping PAROLE/SIMPLE lexicons onto Lemon/LexInfo involves three tasks. Firstly, the original PAROLE/SIMPLE model expressed in the DTD needs to be mapped onto the Lemon model. This can be seen as the lexicon format mapping. Secondly, all descriptive elements defined by PAROLE/SIMPLE lexicons are mapped onto the LexInfo ontology. This includes language dependent descriptive elements and common elements⁶. This broadly corresponds to the ontology mapping part. Finally, lexical entries are mapped.

The resulting dataset is organized into three files. One contains the PAROLE/SIMPLE Ontology which essentially imports Lemon and LexInfo ontologies and adds 'PAROLE elements' (classes and/or properties) whenever these could not be mapped. The other two files collect the Spanish and Catalan lexical entries.

⁵ Especially GENELEX, PAROLE and SIMPLE.

⁶ Note that, whereas PAROLE lexicons are structurally compatible, in certain aspects they are semantically idiosyncratic as each lexicon defines its own 'descriptive' elements. Thus for example, subcategorization frames are defined in each lexicon without any reference or relation to the others. In contrast, SIMPLE lexicons go one step further and define a set of shared descriptive elements.

In the following lines we describe the clues of the mapping process and highlight some of the benefits obtained.

2. From PAROLE/SIMPLE model to Lemon

The strategy followed when mapping PAROLE/SIMPLE model onto Lemon can be summarized as follows:

Elements from the DTD were mapped onto Classes. Whenever possible, Lemon (and LexInfo) classes were used. Otherwise, new classes were created. For example: PAROLE *Description* elements become *lemon:Frames*. In contrast, the *parole:Connotation* class was created as a subclass of *parole:Element* and *lemon:PropertyValue* as shown in Figure 1. Note that many PAROLE/SIMPLE elements are not mapped and simply disappear in the target model. This is partially due to the fact that RDF allows a better modeling and they are no longer needed.

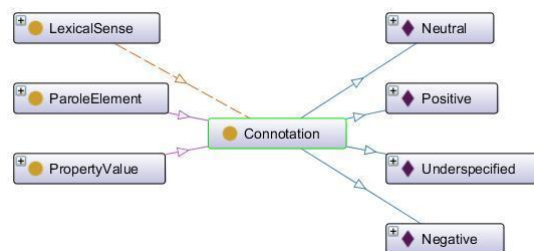


Figure 1 'Adding Classes'

Attributes from the DTD were mapped onto Properties. Again, whenever possible, Lemon or LexInfo properties were used. For example: PAROLE *MuS/@gramcat*⁷ becomes *lexinfo:partOfSpeech*.

Values. When the PAROLE/SIMPLE DTD establishes the set of values for a given attribute, these values are mapped onto the corresponding LexInfo values. For example: the PAROLE pair: "NOUN" + "COMMON" simply translates as *lexinfo:commonNoun* as shown in Figure 2.

Parent/child relations between elements in the DTD were mapped onto relevant Properties. For example: the parent/child relation between a PAROLE verbal *Construction* and its subject *InstantiatedPositionC* element becomes *lexinfo:subject property*.

⁷ We use XPath expressions when referring to source data.

```

<MuS
  id = "0001"
  gramcat = "NOUN"
  gramsubcat = " COMMON"
  ...
</MuS

```

"0001" lexinfo:partOfSpeech lexinfo:commonNoun;

Figure 2 'Attribute mapping'

IDREFs pointing mechanisms between elements in the DTD became properties. For example: the relation between PAROLE morphological and syntactic units (*MuS* & *SynUs*) is expressed by means of the *lemon:synBehaviour* property as shown in Figure 3.

```

<!-- Morphological Units-->
<!-- Syntactic Units-->
<MuS
  id= " 001"
  gramcat = "VERB"
  gramsubcat = " MAIN"
  synulist = "SynU-001"
  ...
</MuS>
<SynU
  id= "SynU-001"
  ...
</SynU>

```

lemon:synBehaviour →

Figure 3 'Mapping the IDREF pointing mechanism'

Though the mapping process implied a considerable effort we think the task was worth it. The source model (DTD) and common descriptive elements are already mapped and can be reused by other languages. Lexical entries and language dependent data in source lexicons will require additional mapping processes. However, this task can benefit from already defined conversion templates. The conversion templates defined in this task can be reused when mapping lexical entries from different languages and sources. Figure 4 shows part of the XSL template used to map PAROLE features to LexInfo ontology.

3. Some benefits: syntax/semantic linking

Lemon model simplifies the original PAROLE/SIMPLE model in a good number of aspects. This is partly due to the use of RDF which allows for a more compact and efficient representation. The case of syntax/semantic mappings is particularly interesting. The original PAROLE/SIMPLE data include a complex machinery to define syntactic subcategorization frames and semantic argument struc-

tures. In the former case, we have to deal with a large set of related elements: *SynU*, *Description*, *Construction*, *Self*, *InstantiatedPositionC*, *PositionC*, *SyntaxmaNT*, etc. The relation among these elements is established by means of the parent/child relation mechanism or ID/IDREF pointing mechanism as exemplified in Figure 5.

```

<xsl:when test="./@value = 'VERB'">
  <xsl:text> lexinfo:partOfSpeech lexinfo:VerbPOS </xsl:text>
</xsl:when>
<xsl:when test="./@value = 'PASSIVE'">
  <xsl:text> lexinfo:voice lexinfo:passiveVoice </xsl:text>
</xsl:when>
<xsl:when test="./@value = '1'">
  <xsl:text> lexinfo:person lexinfo:firstPerson </xsl:text>
</xsl:when>
<xsl:when test="./@value = 'GCOMMON'">
  <xsl:text> lexinfo:gender lexinfo:commonGender </xsl:text>
</xsl:when>
<xsl:when test="./@value = 'SINGULAR'">
  <xsl:text> lexinfo:number lexinfo:singular </xsl:text>
</xsl:when>

```

Figure 4 'Mapping features'

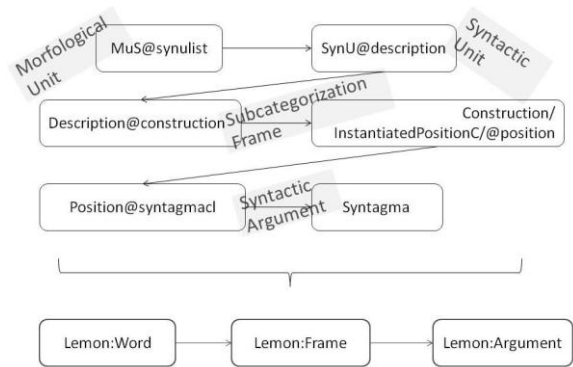


Figure 5 'Subcategorization information'

Similarly, argument structure representation is also complex and, again, we find a good number of elements involved: *PredicativeRepresentation*, *Predicate*, *Argument*, *InfArg*, *SemanticRole*, etc.

Syntax semantic linking in the PAROLE/SIMPLE model is even more complex and, in most cases, useless. Syntactic frame descriptions and semantic predicate descriptions are completely separated. The former involve syntactic arguments whereas the latter involve semantic arguments with no relation at all between them. Syntax/semantic relations are expressed by means of additional elements: the *Correspondence* element and its 'descendants'. *Correspondence* are global elements that point to *SimpleCorrespondArgPos* elements which are the eventual holders

of the syn/sem argument linking. Since *SimpleCorrespArgPos* elements are global, the linking is defined not in terms of arguments IDs but in terms of the position they occupy in the syntactic frame and the semantic predicate. Note in addition (see Figure 6) that neither the syntactic frame nor the predicate involved are at hand.

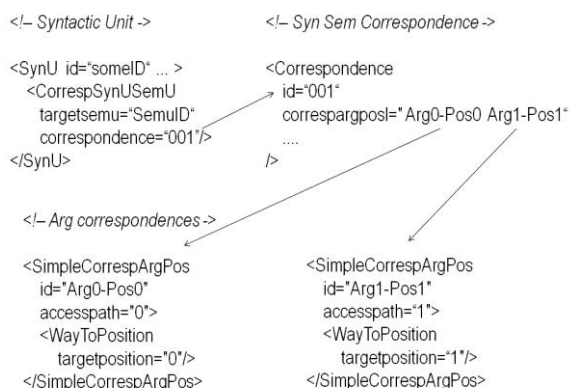


Figure 6 ‘Syn / sem linking in PAROLE/SIMPLE’

Lemon model allows defining all these things in a much easier way, essentially:

Description, Construction & Self elements are mapped to *lemon:Frame* class and related onto the relevant entry by means of the *lemon:synBehaviour* property.

InstantiatedPositionC, Position & Syntagmas are mapped onto *lemon:Argument* class and related to the relevant *lemon:Frame* via some *lemon:synArg* relation.

PredicativeRepresentation & Predicate are also mapped onto *lemon:Frame*

Argument, SemanticRole & InfArg become *lemon:Argument* class and link to relevant *lemon:Frame* via some *lemon:semArg* relation.

A simplified entry for the English verb ‘write’ can be found in Figure 7. Figure 8 gives a partial graphical representation. There we can see that both the syntactic frame and the lexical sense point to ARG0 and ARG1 instances. In the former case, the frame links to its arguments by means of *subject* and *object* properties. In the latter case, the lexical sense links to its arguments by means of *agent* and *patient* properties. Finally, arguments are also specified for a semantic template (Human & SemioticArtifact respectively) and syntactic realization (NP in both cases).

```

## Lexical Entry
lex:write a lemon:LexicalEntry ;
lexinfo:partOfSpeech lexinfo:mainVerb ;
lemon:synBehaviour lex:write_transitive;
lemon:sense lex:write_SymbolicCreation .

## Lexical Senses
lex:write_SymbolicCreation a lemon:LexicalSense ;
parole:template parole:SymbolicCreation ;
parole:roleAgent lex:write_ARG0 ;
parole:rolePatient lex:write_ARG1 .

## Frames
lex:write_transitive rdf:type owl:Thing ;
rdf:type lex:Transitive ;
lexinfo:subject lex:write_ARG0;
lexinfo:directObject lex:write_ARG1 .

## Argument info
lex:write_ARG0 a lexinfo:Subject ;
lemon:constituent lex:NP
parole:template parole:Human .
lex:write_ARG1 a lexinfo:DirectObject ;
lemon:constituent lex:NP
parole:template parole:SemioticArtifact .

```

Figure 7 ‘A simplified entry for the English verb write’

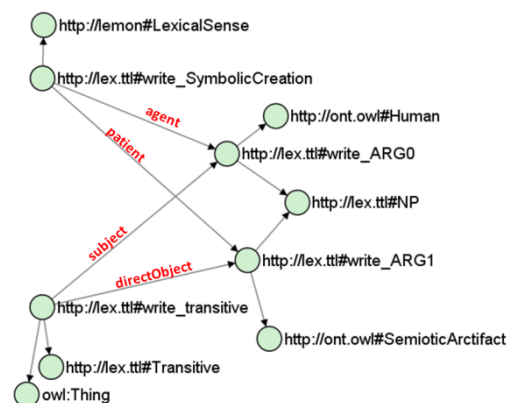


Figure 8 ‘Simplified Syn/Sem linking’

4. Some benefits: subcategorization frames

Each original PAROLE lexicon defines the set of subcategorization frames for a particular language. Contrary to semantic descriptions, syntactic descriptions are essentially language dependent. Thus whereas all lexicons share the same set of semantic descriptive elements (domain, semantic class, semantic relations, etc) such homogeneity was not defined in the syntactic layer. This means that subcategorization information cannot be easily shared among the lexicons.

Basically, this is due to the fact that PAROLE aimed at being a flexible model to accommodate different approaches. This is welcome but proves problematic when addressing interoperability among resources. LexInfo defines a subcategorization ontology based on the Lemon model. As we saw, Lemon includes the notion of *Frame*. Frames are indicated with the *synBehaviour* property and their arguments with the property *synArg*. LexInfo defines subproperties of *synArg* to represent the syntactic functions of arguments and organizes frames into subclasses. Our mapping to LexInfo implied mapping PAROLE subcategorization frames onto this model (*Description* elements and their ‘descendants’). The mapping process was done in two steps. First, we defined a style sheet converter that reads our PAROLE XML lexicon and for each *Description* element it generates a new *Frame*. Consequently, all newly created frames were treated as subclasses of the general *lemon:Frame*. Second, we collapsed some frames into one single class⁸, thus simplifying the model, and organized them in the LexInfo ontology. As a result, the PAROLE ontology becomes lighter than the original model and allows queries that were otherwise impossible in the original PAROLE lexicon; for instance we can easily get all ‘control’ verbs; verbs with a sentential complement; verbs with an indirect object, etc.

5. Some benefits: exploitation

The most difficult problem of the original PAROLE/SIMPLE lexicons is exploitation and management. When moving from the original PAROLE/SIMPLE model to a relational database, we end up with a complex database with a huge number of related tables⁹. Having PAROLE/SIMPLE lexicons in a database means managing lots of tables and very often we need to split complex queries into several sub queries [6]. Note, for example, that getting the senses of a given lemma is not easy and we need a complex SQL query involving up to six different tables. Similarly, a query such as “find the lemma and template of all senses with a negative connotation” is a real challenge in the original PAROLE/SIMPLE lexica. Such a query is quite simple

⁸ For example, the original Spanish lexicon includes 12 intransitive prepositional *Descriptions*, one for each bounded preposition. All these frames are mapped to IntransitivePP Frame as the information about the preposition is encoded by means of a property attached to the PP argument.

⁹ Our PAROLE/SIMPLE database included 223 tables.

in RDF as shown in Figure 9. The results are given in Figure 10.

```
SELECT ?form ?template WHERE {
  ?entry lemon:form [ lemon:writtenRep ?form ] .
  ?entry lemon:sense ?sense .
  ?sense parole:connotation parole:Negative .
  parole:template ?template .
}
```

Figure 9 ‘SPARQL sample query’

Results	
form	template
coco	◆ parole:TemplLivingentity
finolis	◆ parole:TemplHuman
asno	◆ parole:TemplHuman
animal	◆ parole:TemplHuman
animal	◆ parole:TemplHuman
gentuza	◆ parole:TemplHumanGroup
pomposidad	◆ parole:TemplQuality
escasez	◆ parole:TemplAmount
escondrijo	◆ parole:TemplLocation
promiscuidad	◆ parole:TemplQuality
cabrón	◆ parole:TemplHuman
escabrosidad	◆ parole:TemplQuality
tullido	◆ parole:TemplHuman
desmemoriado	◆ parole:TemplHuman
estrechez	◆ parole:TemplQuality
alarmar	◆ parole:TemplExperienceEvent
inconsistencia	◆ parole:TemplQuality
pereza	◆ parole:TemplPsychproperty
escepticismo	◆ parole:TemplPsychproperty
hastiar	◆ parole:TemplExperienceEvent
monstruo	◆ parole:TemplHuman
monstruo	◆ parole:TemplLivingentity
fracasado	◆ parole:TemplHuman
drama	◆ parole:TemplAbstractEntity

Figure 10

6. The sources

The Ontology and both the Spanish and Catalan lexicons are distributed under CC Attribution 3.0 Unported license. These datasets are published in the Data Hub (<http://datahub.io/dataset/parole-simple-ont>) and can be downloaded in both XML RDF format and RDF Turtle format.

The Spanish lexicon contains 199,466 triples with 7,572 lexical entries fully annotated with syntactic and semantic information distributed as follows: 5,659 common nouns, 729 proper nouns, 859 adjectives and 325 verbs. The lexicon contains 11,430 LexicalSenses.

The Catalan lexicon contains 343,714 triples and 20,545 lexical entries annotated with syntactic information half of which are also annotated with semantic information. Lexical entries include 3,064 verbs, 13,206 common nouns, 247 proper nouns, 3,101 adjectives

tives and 511 adverbs. The rest belong to closed categories. The lexicon contains 11,813 LexicalSenses.

Table 1 lists the properties assigned to LexicalSenses¹⁰ in both lexicons.

Property	Spanish	Catalan
id	11430	11813
template	9924	10782
example	9727	10443
semanticClass	8987	10027
semanticRelation	15808	23835
countability	6827	5573
semanticFeature	3222	4328
semanticRole	2294	4381
copulaType	971	
connotation	1314	1364
adjType	979	715
comment	1506	8388
domain	107	56
gradable	246	
definition		10658
TOTAL	73342	102363

Table 1 ‘Triples assigned to LexicalSense’

7. Summary and conclusions

The dataset described here is the result of mapping PAROLE/SIMPLE Spanish and Catalan lexicons onto Lemon model following the LexInfo ontology. The mapping implied three main tasks: the lexicon format mapping (from DTD to Lemon model), the ontology mapping (from ‘descriptive’ elements to LexInfo ontology) and the mapping of lexical entries.

This work may help and encourage other PAROLE/SIMPLE lexicons to take the same way. The Lemon version of PAROLE model (DTD) is already mapped and all shared descriptive elements are integrated with LexInfo ontology. Everything can be reused by other languages. In addition, new lexicons can benefit from conversion templates and only need to address language particular descriptions. Linked Open Data is the natural scenario for a multilingual resource such as the PAROLE/SIMPLE lexicons.

The resulting lexicons benefit from standardization and Linked Open Data principles. The fact that source data categories are mapped onto the LexInfo ontology which in turn is linked to ISocat¹¹ is a step forward in standardization and interoperability.

From our experience we conclude that XML (essentially DTDs) is not well suited for modeling purposes as it allows for a number of syntactic alterna-

tives and conveys semantic ambiguity. In addition, XML proves inefficient when relating resources. This is crucial in a scenario where references to external resources are essential to guarantee interoperability. RDF overcomes some of the problems met with XML. The use of RDF (especially URIs) proves essential for datasets such as original PAROLE/SIMPLE lexicons where interoperability was crucial. The resulting data is lighter and better suited for exploitation. In addition, it eases further extensions and links with external resources such as WordNet, lemonUby, DBpedia etc.

8. Acknowledgements

The resources reported in this paper were developed with the support of METANET4U: Enhancing the European Linguistic Infrastructure, (2011-2013), funded by UNER - Competitiveness and Innovation Framework Program, (CIP-PSP-270893).

We thank the Institut d’Estudis Catalans and the GilcUB from the University of Barcelona as the creators of the original Catalan and Spanish lexicons.

References

- [1] Marie-Hélène Antoni-Lay, Gil Francopoulo, Laurence Zaysser. A generic model for reusable lexicons: the GENELEX project. *Literary and linguistic computing*. Vol 9. Num 1. Pages 47-54. (1994).
- [2] Paul Buitelaar, Philipp Cimiano, Peter Haase, Michael Sintek: Towards Linguistically Grounded Ontologies. *ESWC 2009*: 111-125
- [3] Philipp Cimiano, Paul Buitelaar, John McCrae, Michael Sintek: LexInfo: A declarative model for the lexicon-ontology interface. *J. Web Sem.* 9(1): 29-51 (2011)
- [4] Philipp Cimiano, Paul Buitelaar, John McCrae., & Michael Sintek, M. (2011). *LexInfo: A Declarative Model for the Lexicon-Ontology Interface*. *Web Semantics: Science, Services and Agents On The World Wide Web*, 9(1).
- [5] Gil Francopoulo, Núria Bel, Monte George, Nicoletta Calzolari, Monica Monachini, M Pet, Claudia Soria. 2007 *Lexical Markup Framework: ISO standard for semantic information in NLP lexicons*. GLDV (Gesellschaft für linguistische Datenverarbeitung), Tübingen
- [6] Marta Villegas, Núria Bel. (2002) .From DTD to relational dB. An automatic generation of a lexicographical station out off ISLE guidelines. In *Proceedings LREC*, 2002.
- [7] Sanfilippo et al., 1996 Preliminary Recommendations on Sub-categorization <http://www.ilc.cnr.it/EAGLES96/synlex/synlex.html>
- [8] McCrae J., D. Spohr, P. Cimiano (2011) *Linking Lexical Resources and Ontologies on the Semantic Web with Lemon* *Proceedings of the 8th European Semantic Web Conference, Lecture Notes in Computer Science*, Springer, Volume 6643, pp.245-259.

¹⁰ Semantic relations and semantic roles are grouped. The object of ‘semantic relation’ triples is always another LexicalSense.

¹¹ <http://www.isocat.org/>