

Ressources interoperability

EXPLOITING LEXICOGRAPHIC DATA TO AUTOMATICALLY GENERATE DICTIONARY EXAMPLES

María José Domínguez / Miguel Anxo Solla / Carlos Valcárcel



Fundación **BBVA**

Content

1

Core principles

The MultiGenera and Multicomb projects

2

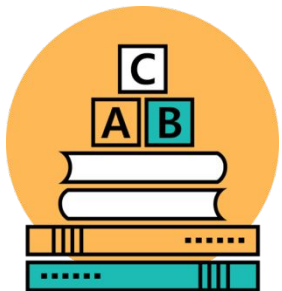
PORTLEX

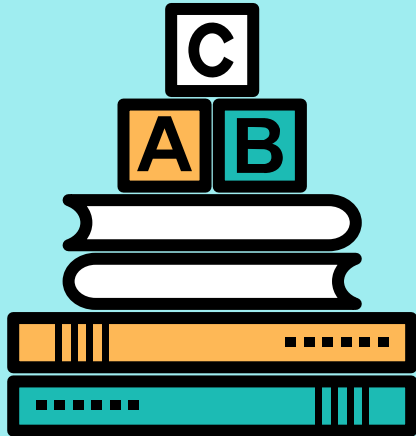
The main lexicographic source

3

Methodology & tools

Functionalities and uses of the APIs, LEMATIZA, COMBINA and XERA





1. Core principles

The MultiGenera and MultiComb projects

Two **linked** projects



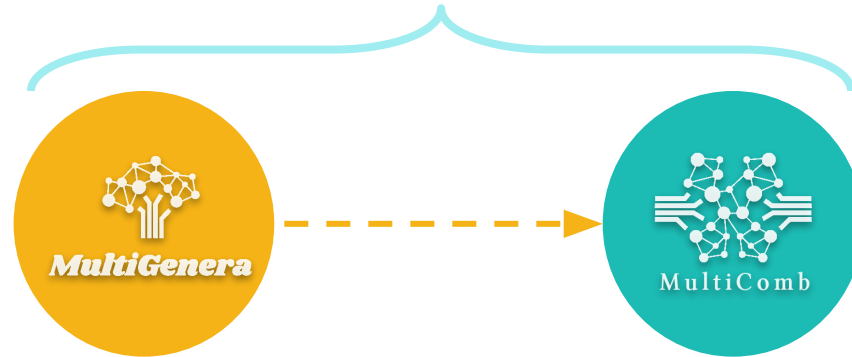
ES



DE



FR



MultiGenera

Nominal phrase generation

MultiComb

Sentence context generation

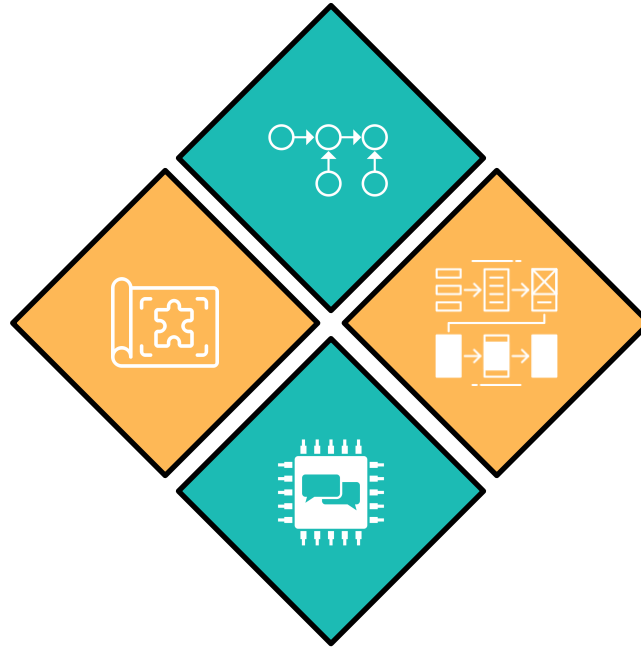
Approaches applied

Valency Grammar

Arguments
Surface realizations
Semantic roles
Combinatory patterns

Meaning to Text Theory

Lexical functions
Semantic representation
Government patterns



Prototypes Theory

Lexical prototype
Semantic prototype

Natural Language Processing

Natural Language Generation
Information extraction
Lexical mining

Combined method



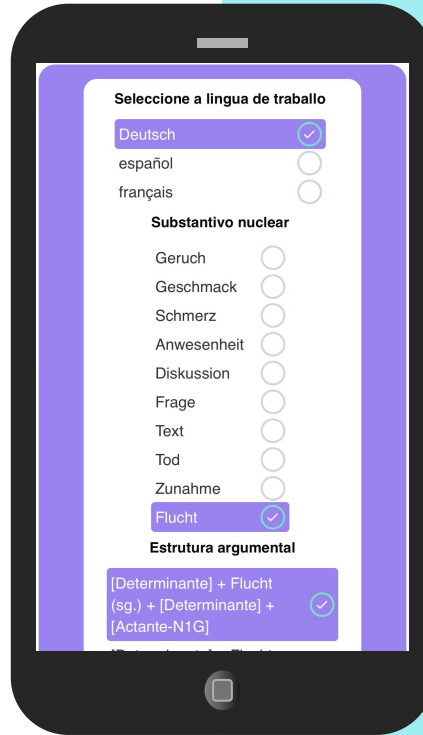
Analysis



Extraction



Evaluation



Ressource interoperability



External

PORTLEX
WordNet
FreeLing



Internal

APIs
Lematiza
Combina

What can it be useful for?



1.
For building customized examples.

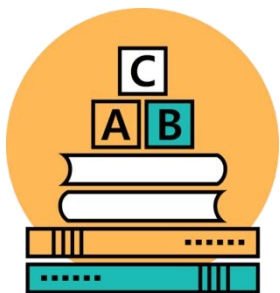


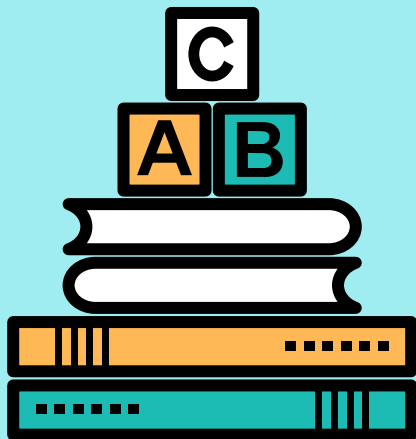
3.
For building lexical extraction tools.

2.
For creating reusable lexical packages.



4.
For studying restrictions within the noun phrase.



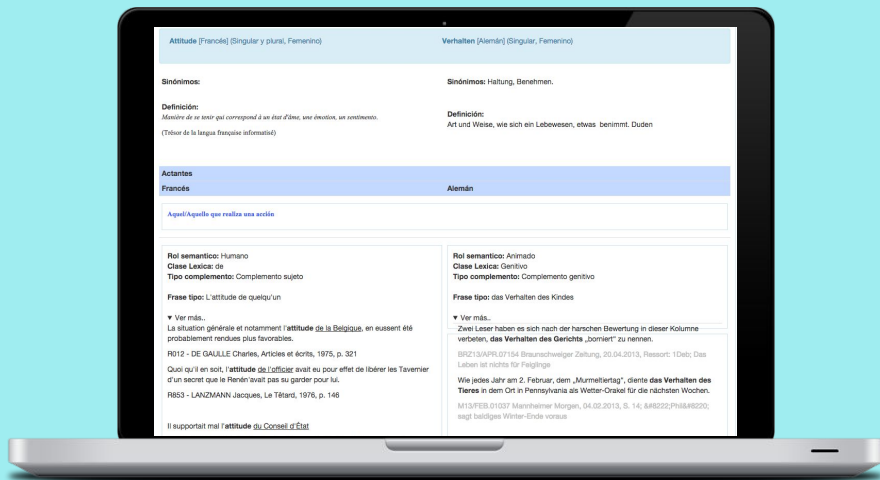


2. PORTLEX

Our main lexicographic source



A valency dictionary of the noun phrase



Valency



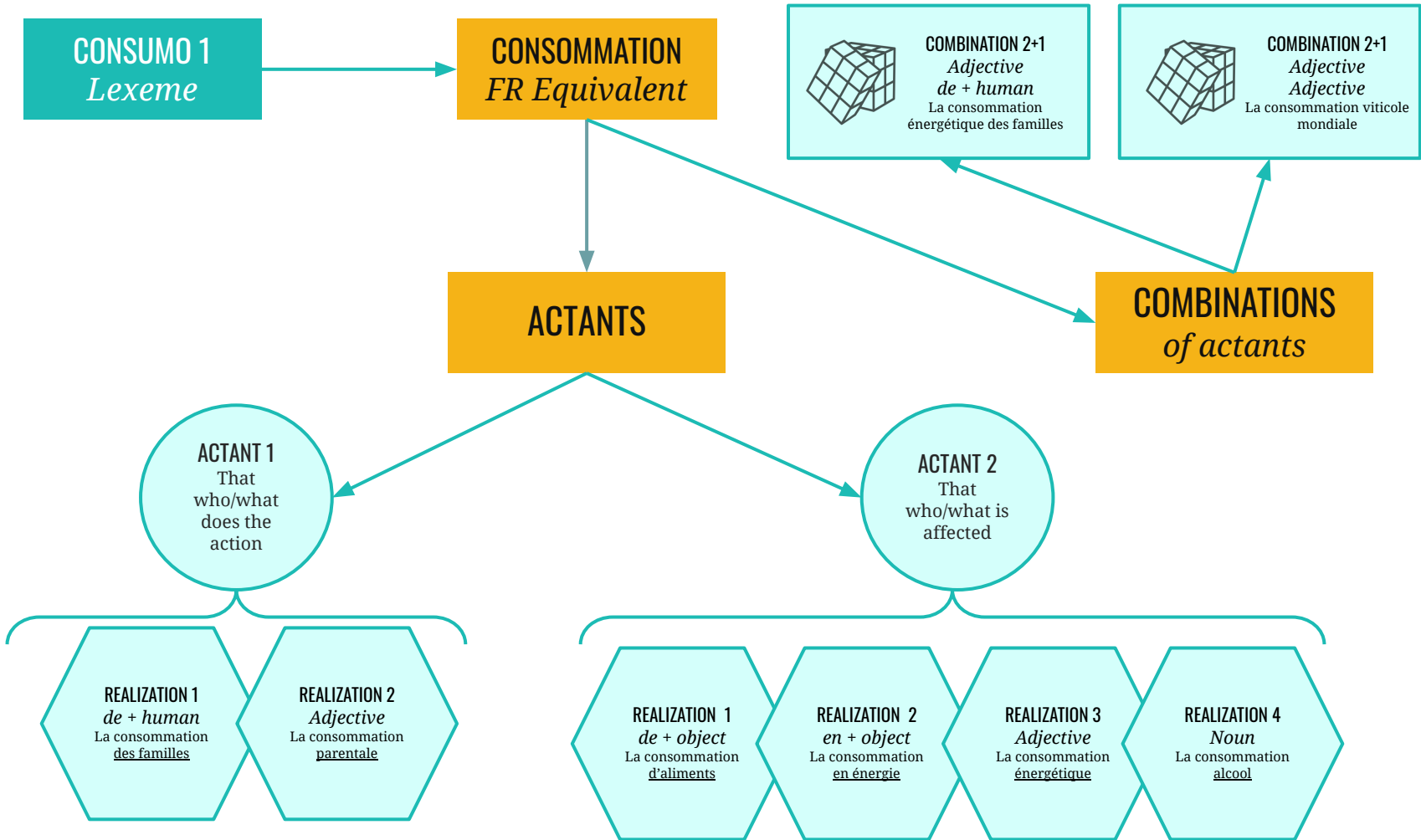
Collaborative



Cross-lingual



Modular





Followers

They follow PORTLEX updates on the web and/or on social media. They can interact with users and editors, but they are not registered as users since they don't consult the dictionary.



Editors

Specialists in a particular language with training in valency grammar. They edit entries in the dictionary database..



Reviewers

Specialists in a particular language with training in valency grammar. They validate the data entered by the editors.



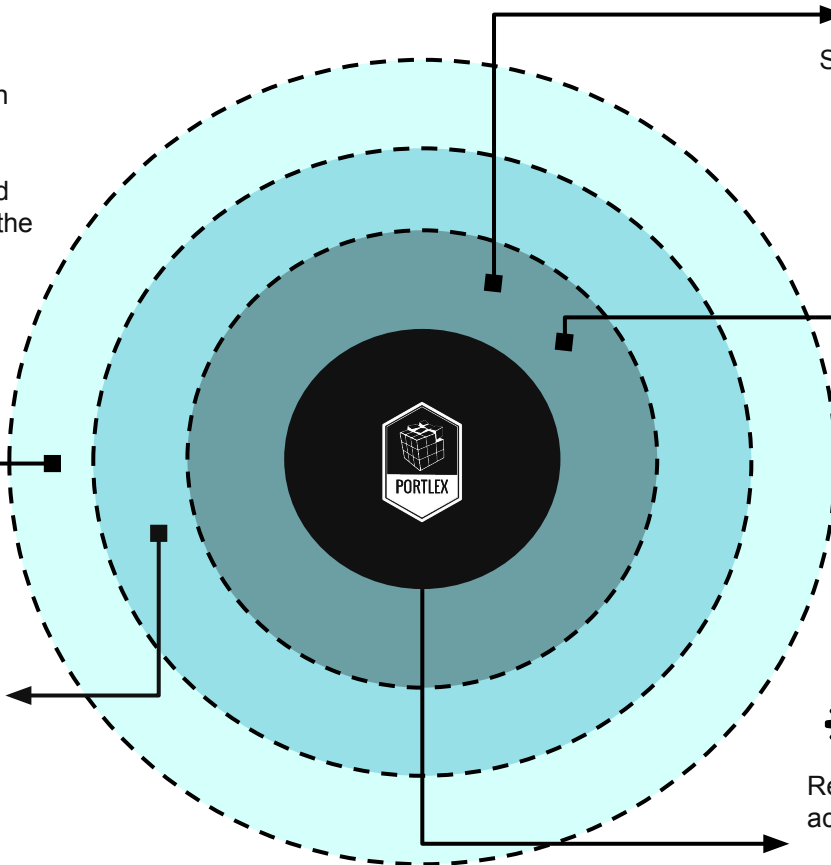
Administrators

Researchers who manage the database, access and user roles.

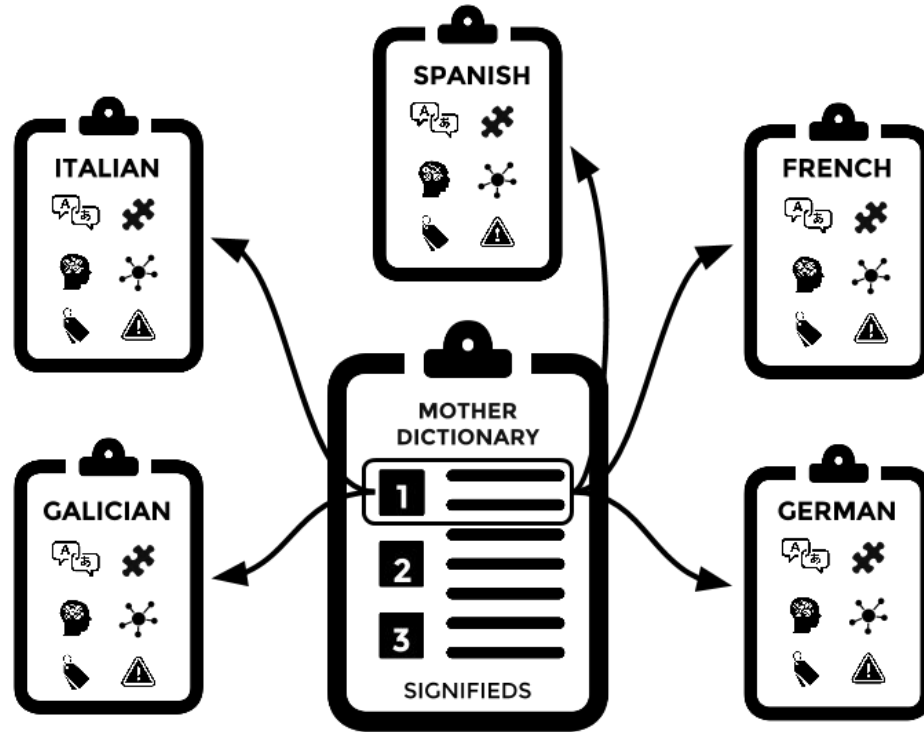


Users

Experts or general users who usually do searches in the dictionary. They can interact with other users and with the editors. They can also become editors after a training process.

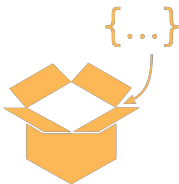


Modular & cross-lingual



Difficulties encountered

when working with corpora



The compilation of surface realizations

It's a very time demanding process.



The description of argument patterns

Dozens of patterns by noun to be described in detail.

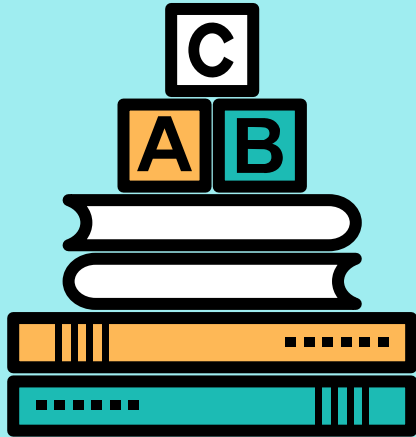


Corpus-extracted data adequacy

Human review is also necessary.



**Semantic
prototypes.**
Could they be
one way out?



3. Methodology & tools

APIs, Lematiza, Combina and Xera.

Combined **method** phases



Lexico-semantic prototyping



Concordances

Pattern description
CQL queries
Slot filling lemmata



Frequencies

Frequent lemmata
Lexical prototypes



Cleansing

Inadequacies
Repetitions



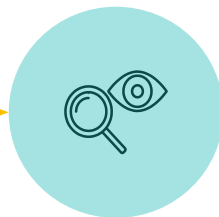
Classification

Semantic
annotation
Semantic
prototyping

An exemple

Det. + Flucht + aus + Noun

Die Flucht aus Sintra



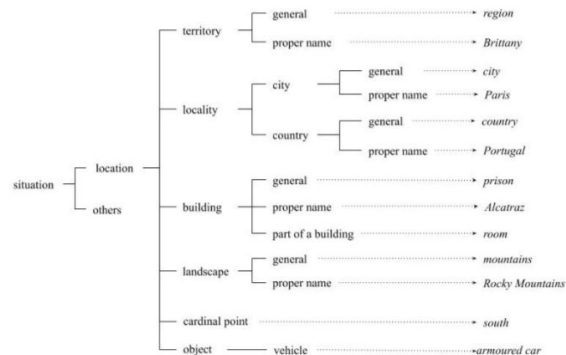
CQL query

DDR
Ghetto
Troja
Haus
Frankreich
Ost-Berlin
etc.

Lexical prototypes

FLUCHT
aus+dative +

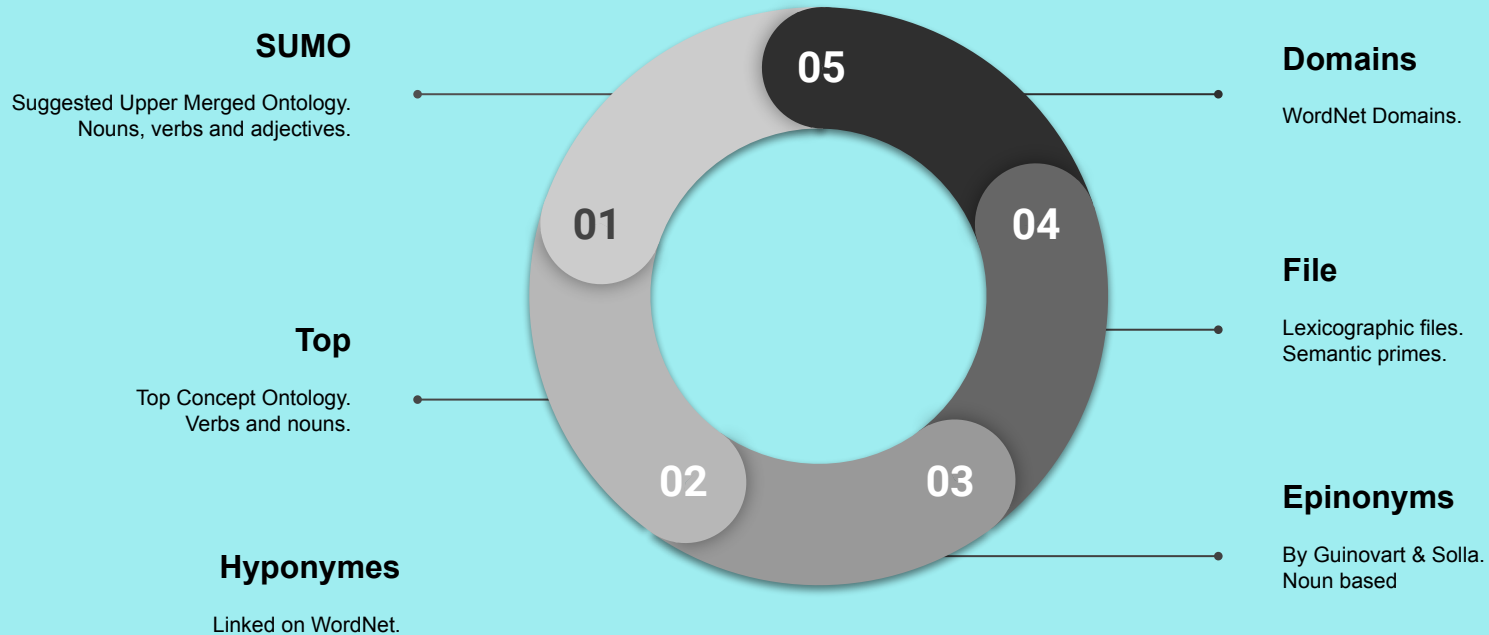
Lexical prototypes	1st Level	2nd Level	3rd Level	4th Level
Warschauer Ghetto	situation	location	territory	
Haus	situation	location	building	
Kriegsgefangenenlager	situation	location	building	
Wohnung	situation	location	building	
Troja	situation	location	locality	proper name
Ost-Berlin	situation	location	locality	proper name
Venedig	situation	location	locality	proper name
Frankreich	situation	location	territory	proper name
Deutschland	situation	location	territory	proper name
Italien	situation	location	territory	proper name



Semantic annotation

Semantic classes

Expansions: resorting to Wordnet



Prototypes expansion



Ontologies

Semantic
matching
LEMATIZA



API queries

Lexical extraction
queries
APIs + COMBINA



Cleansing

Inadequacies
Repetitions



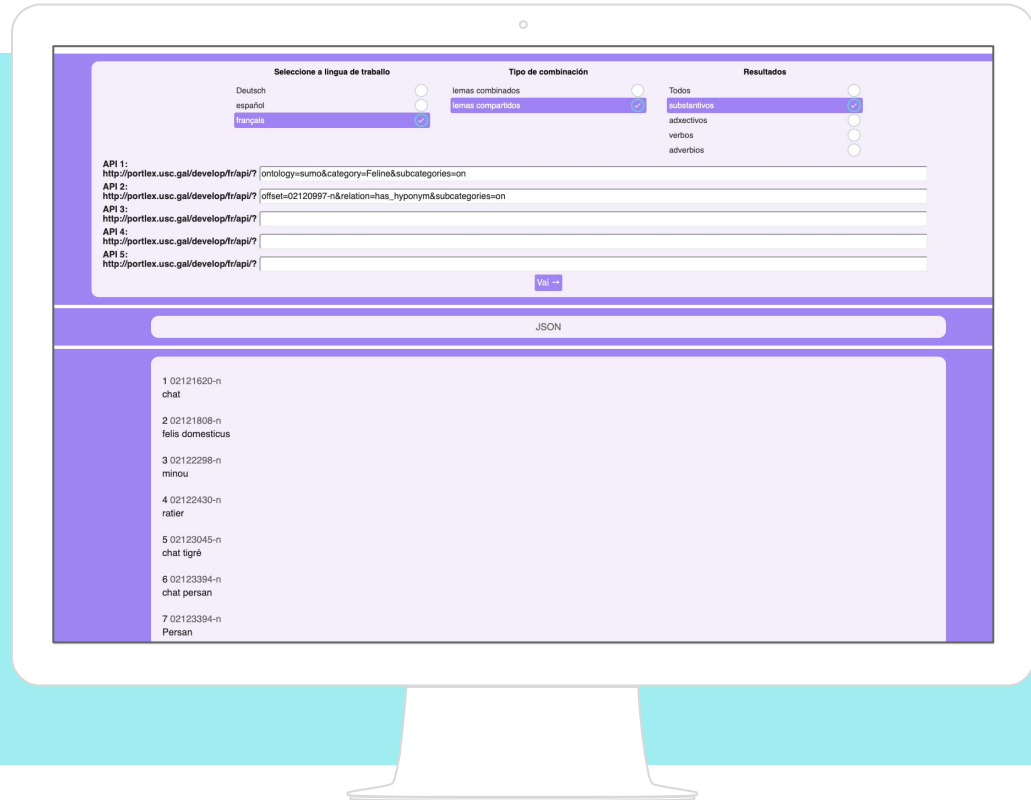
Packaging

Morphological
annotation
Semantic tagging
FLEXIONA

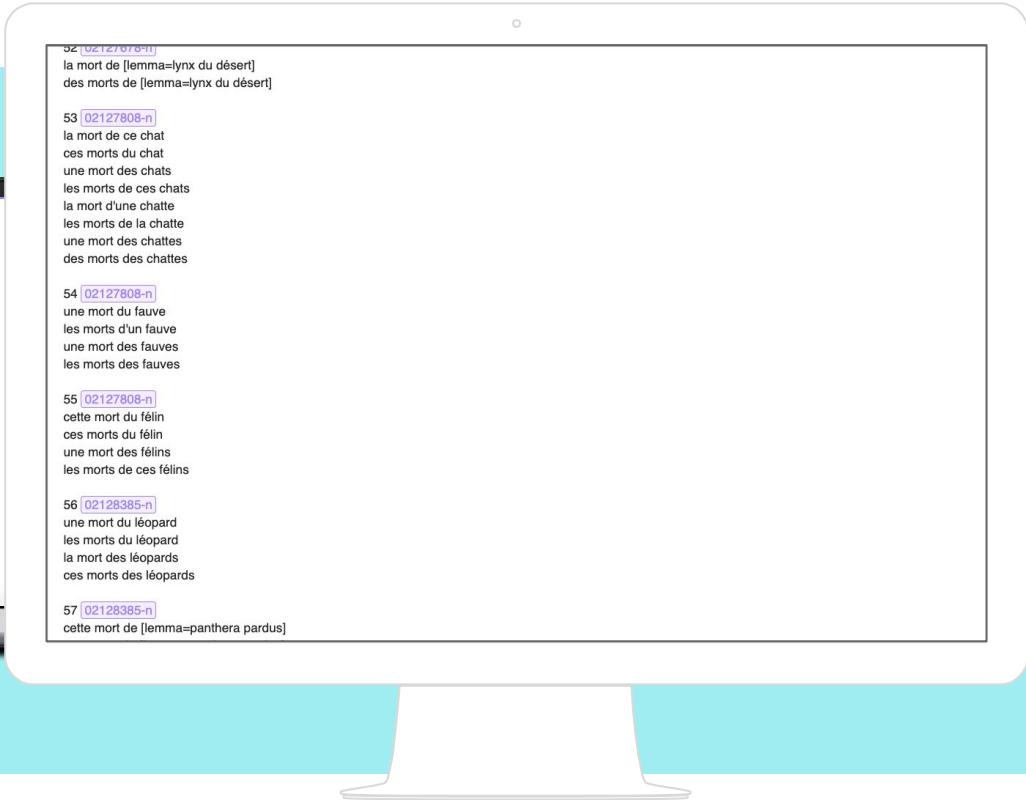
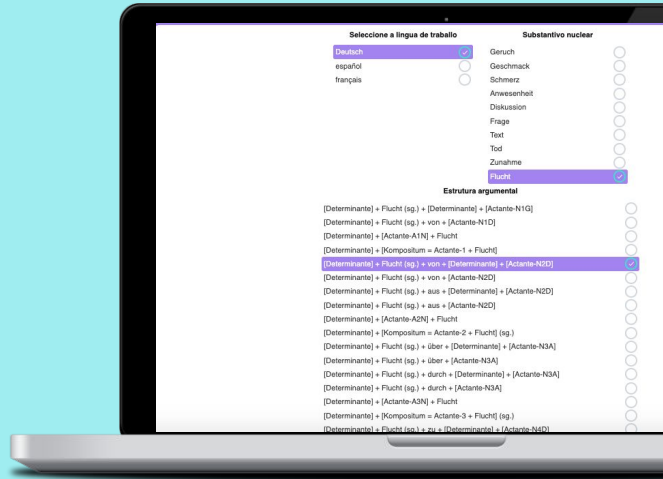
LEMATIZA bit.ly/multigenera-tools

```
1 le odeur de citron
- Actante: citron
- Lema actancial: citron
Offsets:
- 03655838-n an artifact (especially an automobile) that is defective or unsatisfactory
  --- WordNet Domains: factotum + subcategories
  --- SUMO: Artifact + subcategories
  --- Top: Artifact + subcategories | Object + subcategories
  --- Epinonyms: ili-30-00021939-n#artefact + subcategories
  --- Hiperónimo(s): 00021939-n#artefact | artifact Hyponyms + subcategories
  --- Nivel de hiponimia (substantivos e verbos): 5
  --- Ficheiro lexicográfico (substantivos): artifact
- 04966543-n a strong yellow color
  --- WordNet Domains: color + subcategories
  --- SUMO: ColorAttribute + subcategories
  --- Top: Physical + subcategories | Property + subcategories
  --- Epinonyms: ili-30-04916342-n#property + subcategories
  --- Hiperónimo(s): 04965661-n#yellow | yellowness Hyponyms + subcategories
  --- Nivel de hiponimia (substantivos e verbos): 8
  --- Ficheiro lexicográfico (substantivos): attribute
- 05716342-n a distinctive tart flavor characteristic of lemons
  --- WordNet Domains: food + subcategories
```

COMBINA bit.ly/multigenera-tools



XERA bit.ly/multigenera-tools



52 **02127070-f1**
la mort de [lemma=lynx du désert]
des morts de [lemma=lynx du désert]

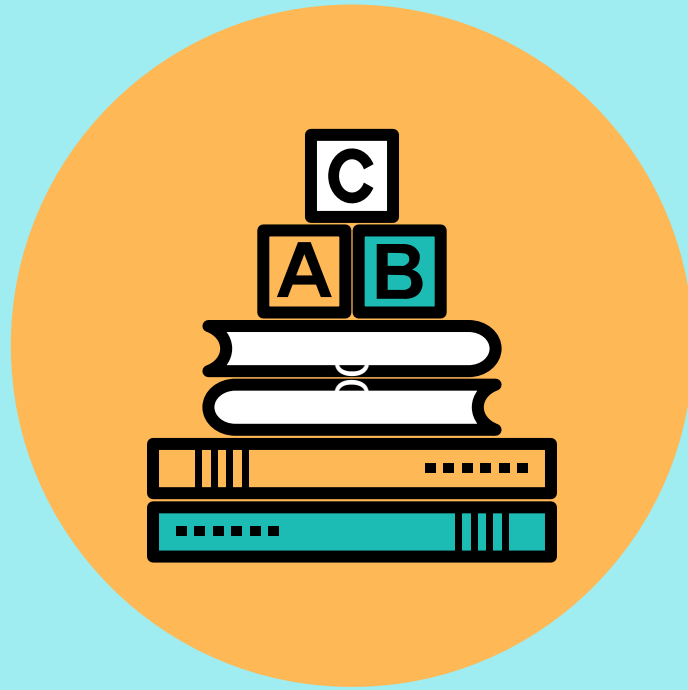
53 **02127808-n**
la mort de ce chat
ces morts du chat
une mort des chats
les morts de ces chats
la mort d'une chatte
les morts de la chatte
une mort des chattes
des morts des chattes

54 **02127808-n**
une mort du fauve
les morts d'un fauve
une mort des fauves
les morts des fauves

55 **02127808-n**
cette mort du félin
ces morts du félin
une mort des félins
les morts de ces félins

56 **02128385-n**
une mort du léopard
les morts du léopard
la mort des léopards
ces morts des léopards

57 **02128385-n**
cette mort de [lemma=panthera pardus]



Thank you! Any questions?

majo.dominguez@usc.es
miguelsolla@uvigo.es
carlos.valcarcel@uvigo.es