



How to use vocabularies to enrich GoTriple ?

Discover
Connect
Collaborate





I am your speaker for this session
JULIEN HOMO



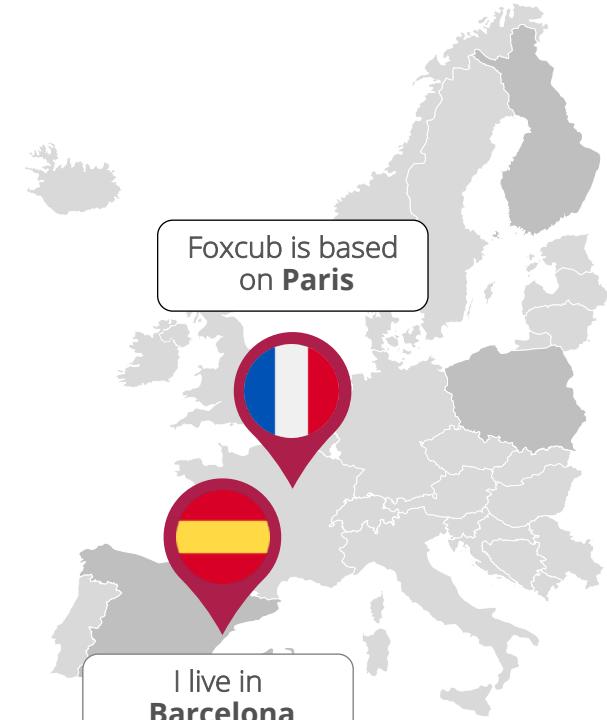
Senior Data Architect
AI Expert



foxcub.
12 years of experience
in Data



Specialized in
Public sector



Discover
Connect
Collaborate



How to use vocabularies to enrich GoTriple ?

- 1 Enriching GoTriple data thanks to vocabularies ?**
- 2 The TRIPLE Vocabulary and MORESS
- 3 The annotation and classification GoTriple services based on vocabularies
- 4 Conclusion & future work
- 5 Questions

Discover
Connect
Collaborate

Definition of a vocabulary in GoTriple



A collection of **semantic resources (or terms)** usually specific to a **subject domain** and created following a **particular methodology**. The semantic resources of a vocabulary are frequently organised hierarchically and may be interlinked forming a whole network of terms.



Thesauri

Vocabularies of terms that are usually domain-specific, organised hierarchically in one or more levels, and containing synonyms, antonyms, etc.

Getty AAT...



Classifications and classification schemes

Presents terms in thematic categories. Classification schemes often contain serial alphanumeric values

DDC, LCC, MORESS...



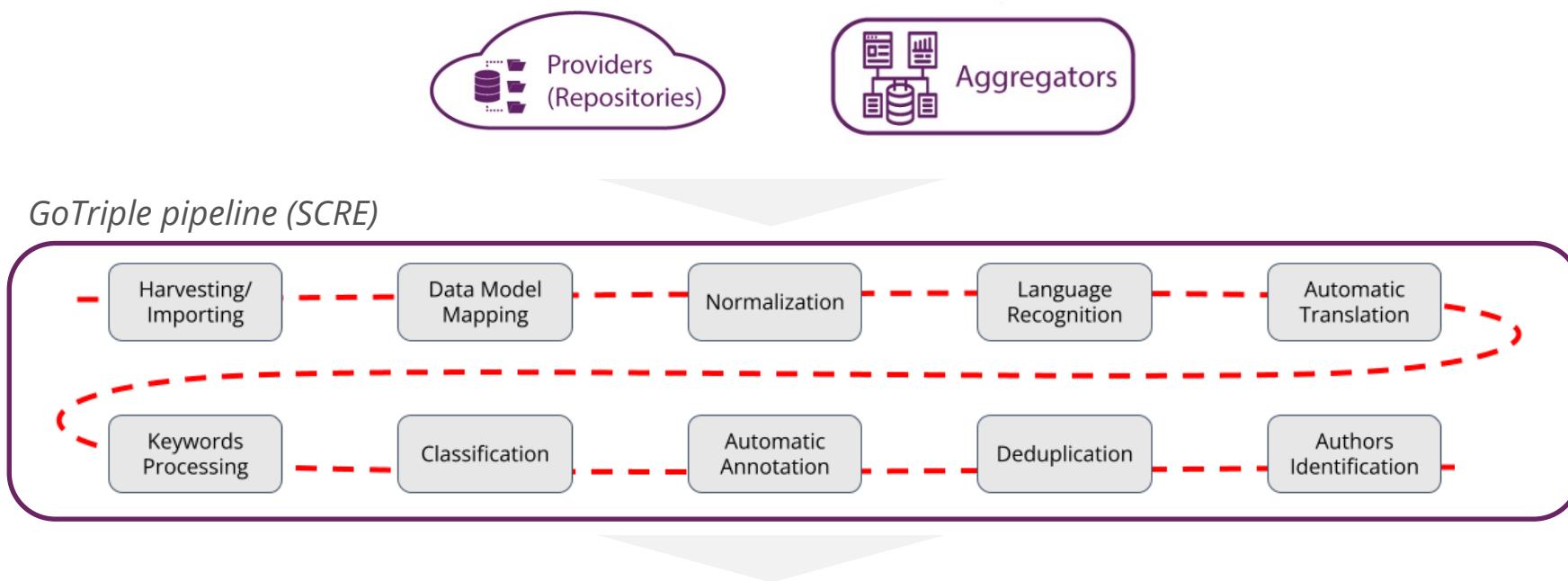
Authority files

List of standardised terms for person names, corporate bodies or concepts commonly used within a domain

VIAF, LCSH, MeSH...

Discover
Connect
Collaborate

The GoTriple platform

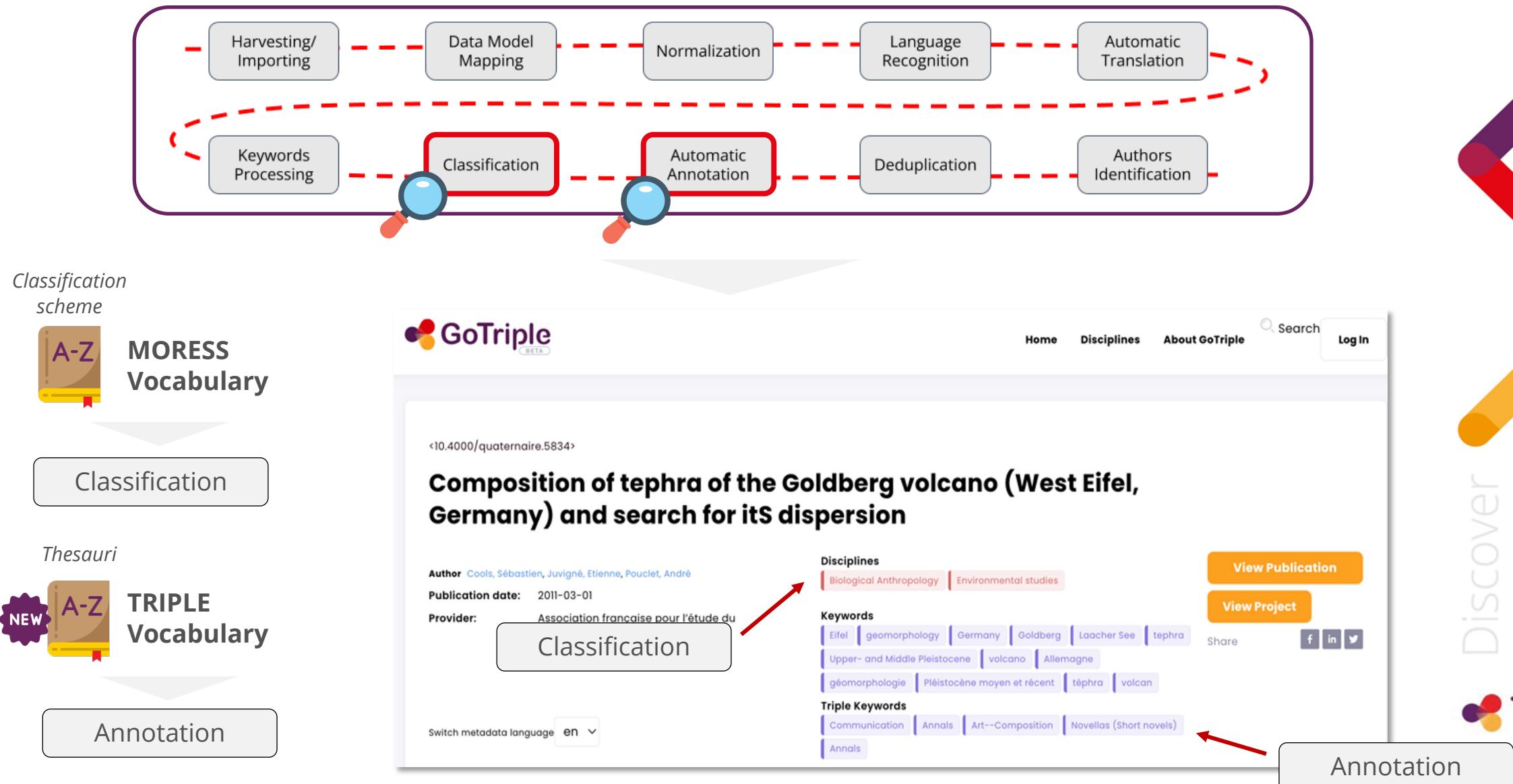


What are the vocabularies involved and where are they (in the GoTriple pipeline) ?

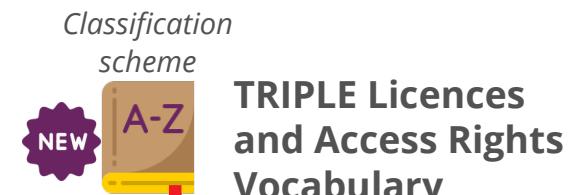
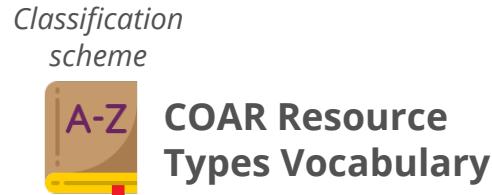
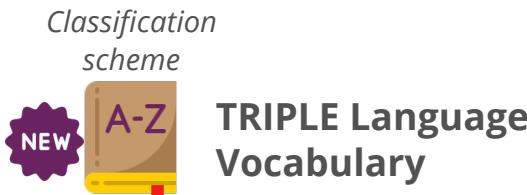
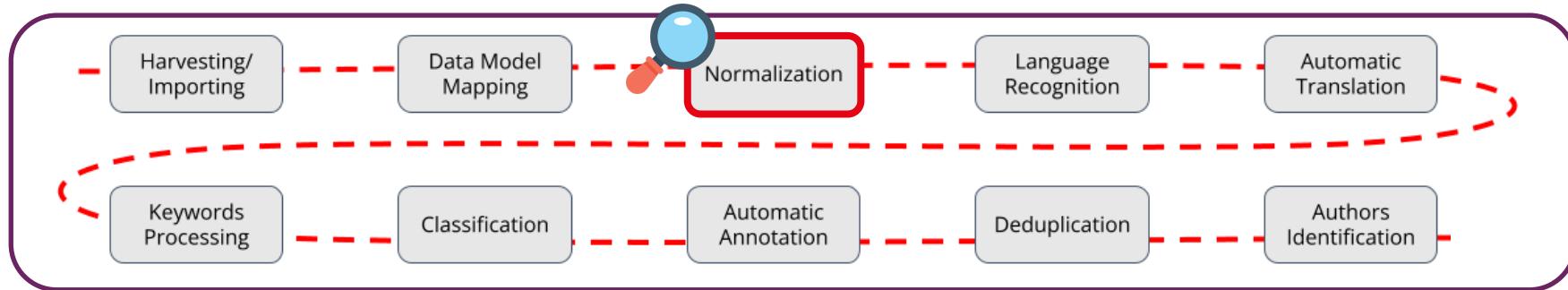
The screenshot shows the GoTriple platform's search interface. The top navigation bar includes Home, Trust Building System, Crowdfunding, Disciplines, Membership, About, Log In, and Sign Up. The search bar displays '750.403 results' for the query 'process'. The search results are presented in a grid format with columns for 'List' and 'Visual'. On the left, there is a sidebar for 'Your Search' with a search input field containing 'process' and a 'Search' button. Below it is a 'Filter Documents' section with a 'Discipline' dropdown set to 'Education' and a 'Search' button. The main results area shows an article titled 'Improvement of the PROCESS OF PROCESS OF PROCESS OF PROCESS OF DISCOUNT' by 'Апостолова Наталья Николаевна' from 'Article, 2018-03-01'. The article abstract is partially visible: 'Расширение составляющих началь в уголовном судопроизводстве и повышение объективности процесса доказывания должно осуществляться в соответствии с логикой, смыслом и духом действующего Уголовно-процессуального законодательства. Ни какие сведения, полученные не в установленном УПК...'.

Discover Connect Collaborate

Classification and annotation: 2 vocabularies



Other vocabularies involved in the GoTriple pipeline



Normalization

Normalization

Normalization

- Croatian (hr)
- Catalan (ca)
- English (en)
- French (fr)
- German (de)
- Greek (el)

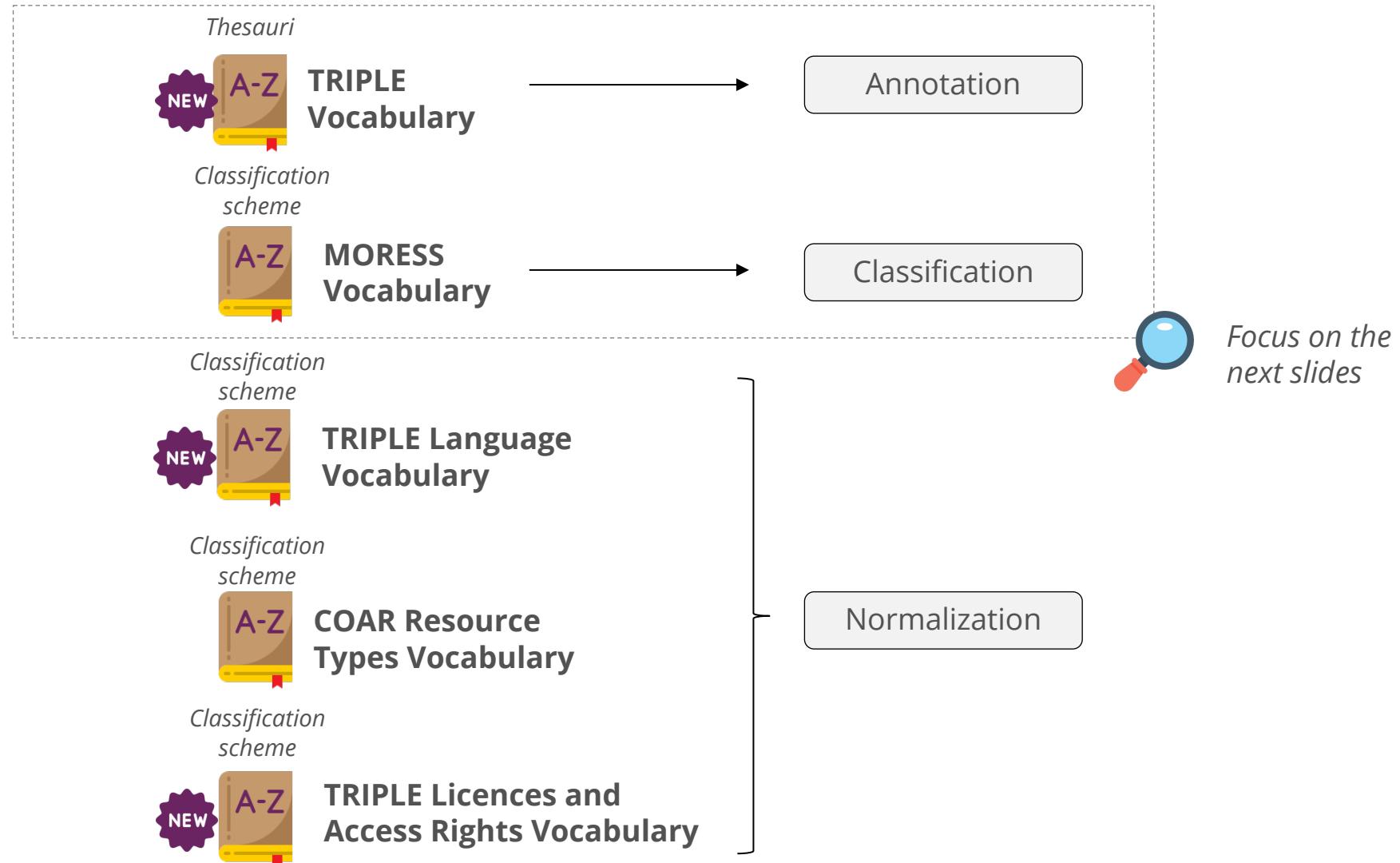
TRIPLE Document Type	Codification	COAR Resource Type
Article	typ_article	Journal Article: https://vocabularies.coar-repositories.org/resource_types/c_6501/
Bibliography	typ_bibliography	Bibliography: https://vocabularies.coar-repositories.org/resource_types/c_86bc/
Blog post	typ_blog-post	Blog post: https://vocabularies.coar-repositories.org/resource_types/c_6947/

TRIPLE License	Codification
CAIRN	lic_cairn
Creative Commons	lic_creative-commons
Various spelling and acronyms, e.g. CCO, CC BY, ..., and full Creative Commons URLs as well.	
Open source	lic_open-source
Various spelling and licence names, e.g. Apache, GPL, BSD, MIT	

Discover Connect Collaborate



There is a strong link between the vocabulary and the service or application that consumes it



How to use vocabularies to enrich GoTriple ?

- 1 Enriching GoTriple data thanks to vocabularies ?
- 2 The TRIPLE Vocabulary and MORESS**
- 3 The annotation and classification GoTriple services based on vocabularies
- 4 Conclusion & future work
- 5 Questions

Discover
Connect
Collaborate

The TRIPLE Vocabulary: overview (1/2)



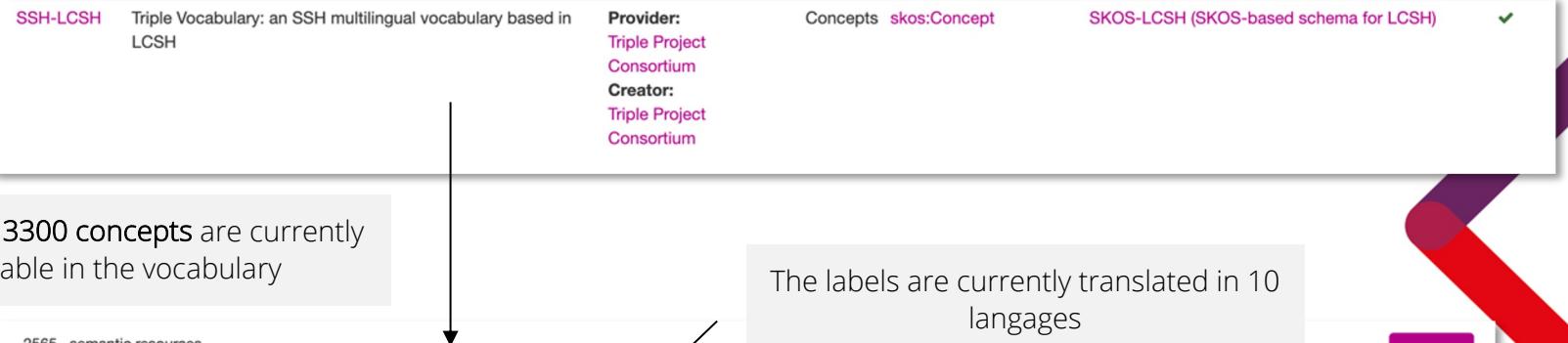
The Triple Vocabulary is published through the Semantics.gr platform

Thesauri



The Triple Vocabulary is a multilingual and hierarchical set of SSH-related concepts. It is a subset of LCSH (Library of Congress Subject Headings) that cover popular SSH aspects and is enhanced with labels in Greek, French, Polish, German, Italian, Portuguese, Spanish, Croatian, Slovenian).

The vocabulary is used for the automatic annotation of the publications hosted in the GoTriple platform



2565 semantic resources

http://semantics.gr/authorities/SSH-LCSH/sh85005581
Anthropology [Flag] Antropología [Flag] antropologija [Flag] Croatian [Flag] Antropologia [Flag] Avθρωπολογία [Flag] Anthropologie [Flag] Anthropologie [Flag] antropologia [Flag] Finnish [Flag]

Close match: [Flag] Anthropology [Flag] ANTHROPOLOGY [Flag] antropologia [Flag] Anthropologie [Flag] Anthropologie [Flag] anthropology [Flag] Antropologia [Flag]

Same as: [Flag] Anthropology [Flag] LCSH

http://semantics.gr/authorities/SSH-LCSH/sh85026423
Civilization [Flag] Civilización [Flag] Civilizacija [Flag] Croatian [Flag] Civiltà [Flag] Πολιτισμός [Flag] Zivilisation [Flag] Civilisation [Flag] sivilisaatio [Flag] beschaving [Flag] Dutch [Flag] Cywilizacja [Flag] Polish [Flag]

Close match: [Flag] Civilisation [Flag] civilization [Flag] civiliatio [Flag] Pays en voie de développement -- Civilisation [Flag] Civilization [Flag] Civiltà [Flag] Zivilisation [Flag]

Same as: [Flag] Civilization [Flag] LCSH

Triple Vocabulary: an SSH multilingual vocabulary based in LCSH

URI: <http://semantics.gr/authorities/vocabularies/SSH-LCSH>

RDF/XML JSON-LD N-triples CSV (Triple)

The vocabulary can be extracted in 4 different formats :
RDF/SKOS, JSON-LD, N-triples and CSV

Category	Concepts
Semantic class	skos:Concept
Provider	Triple Project Consortium
Creator	Triple Project Consortium
Default language	[Flag]
License	Attribution (CC BY 4.0)

```
<rdf:RDF xmlns:dc="http://purl.org/dc/elements/1.1/" xmlns:dcterms="http://purl.org/dc/terms/" xmlns:skos="http://www.w3.org/2004/02/skos/core#">
  <skos:ConceptScheme rdf:about="http://semantics.gr/authorities/vocabularies/SSH-LCSH">
    <dct:title xml:lang="el">Αεξάρνο για ΑΚΕ βασισμένο στο LCSH από την Triple</dctitle>
    <dct:title xml:lang="en">Triple Vocabulary: an SSH multilingual vocabulary based in LCSH</dctitle>
    <dct:creator xml:lang="en">Triple Project Consortium</dct:creator>
    <dct:subject xml:lang="en">Social sciences and humanities</dct:subject>
    <dct:description xml:lang="en">The Triple Vocabulary is a multilingual and hierarchical set of LCSH terms. It is a subset of LCSH (Library of Congress Subject Headings) that cover popular SSH aspects and is enhanced with labels in Greek, French, Polish, German, Italian, Portuguese, Spanish and Croatian. The vocabulary is published as Linked Open Data (LOD) and is available under a Creative Commons Attribution (CC BY 4.0) license. The URIs and their English labels are taken from the 2021 version of the Triple Consortium, using English as the source language.</dct:description>
    <dct:source xml:lang="en">The URIs and their English labels are taken from the 2021 version of the Triple Consortium, using English as the source language.</dct:source>
```

The TRIPLE Vocabulary: overview (2/2)

The data model of the GoTriple Vocabulary is **SKOS**, a standard for data modeling vocabularies and thesauri. SKOS is an RDF Data Model. It defines classes and properties to describe concepts and their relationship. All concepts of the GoTriple Vocabulary are instances of the skos:Concept class.

A **dedicated schema** was created in Semantics.gr based on SKOS to which only the GoTriple Vocabulary conforms. This way, new properties can be easily in the future, even custom ones, without affecting other SKOS vocabularies in Semantics.gr.

<http://semantics.gr/authorities/vocabularies/SSH-LCSH>

The screenshot shows the Semantics.gr platform interface. At the top, there are links for 'Search', 'Vocabularies', 'Derivative Vocabularies', 'Data Models and Schemas', and 'LD & API'. The main content area displays the 'Triple Vocabulary: an SSH multilingual vocabulary based in LCSH' page. It includes a summary box with details like 'Category: Concepts', 'Semantic class: skos:Concept', 'Provider: Triple Project Consortium', 'Creator: Triple Project Consortium', 'Default language: Greek', and 'License: Attribution (CC BY 4.0)'. Below this is a search bar with 'Word or phrase' and 'Search' button, and a link to 'More search options'. A footer section shows '3375 semantic resources' and lists various terms like Anthropologie, Antropologia, Civilizacijā, and Barbarismus, each with multiple language variants (e.g., French, German, Spanish, etc.) and download links for RDF/XML, JSON-LD, N-Triples, and CSV.

Class	skos:Concept			
Properties	qualified name	description	type	multivalued
skos:prefLabel	preferable labels (one per language)	Literal	TRUE (one per language)	
skos:altLabel	alternative labels	Literal	TRUE	
skos:browder	URI Ref to broader concept(s) from the same vocabulary. The property defines the hierarchy. Its symmetric property is skos:narrower.	URI Ref	TRUE	
skos:narrower	URI Ref to narrower concept(s) from the same vocabulary. The property defines the hierarchy. Its symmetric property is skos:browder.	URI Ref	TRUE	
skos:closeMatch	URI Refs to similar resources (with similar meaning) of external vocabularies	URI Ref	TRUE	
skos:exactMatch	URI Refs to exact resources (with exact meaning) of external vocabularies	URI Ref	TRUE	

The MORESS Vocabulary: overview

Classification
scheme



The MORESS vocabulary (Mapping of Research in European Social Sciences and humanities) is mapping to support the development of the European Research Area through improving information on pan-European research in social sciences and humanities. This initiative is part of an EC-funded main program "FP5-HUMAN POTENTIAL - Programme for research, technological development and demonstration on "Improving the human research potential and the socio-economic knowledge base" (1998-2002)"

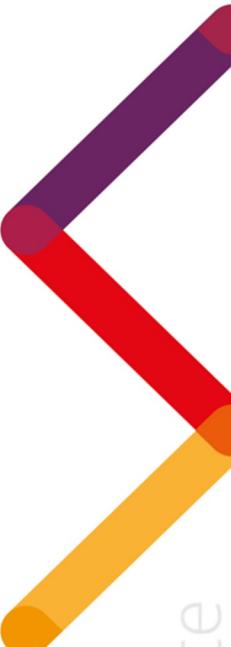
The MORESS Vocabulary **contains 27 categories** which are used to classify the GoTriples document with their associated "**discipline**".

<https://cordis.europa.eu/project/id/HPSE-CT-2002-60060/fr>

MORESS IDs	Labels (FR, EN, ES)	DBpedia / Calenda mapping	HAL mapping
anthro-bio	Biological anthropology	http://dbpedia.org/resource/Biological_anthropology	http://aurehal.archives-ouvertes.fr/subject/shs.anthro-bio
anthro-se	Social Anthropology and ethnology	http://calenda.org/categories.rdf#categorie213 http://dbpedia.org/resource/Social_anthropology	http://aurehal.archives-ouvertes.fr/subject/shs.anthro-se
archeo	Archaeology and Prehistory	http://calenda.org/categories.rdf#categorie293 http://dbpedia.org/resource/Archaeology	http://aurehal.archives-ouvertes.fr/subject/shs.archeo
archi	Architecture, space management	http://dbpedia.org/resource/Architecture	http://aurehal.archives-ouvertes.fr/subject/shs.archi
art	Art and art history	http://calenda.org/categories.rdf#categorie278 http://dbpedia.org/resource/Art_history	http://aurehal.archives-ouvertes.fr/subject/shs.art
class	Classical studies	http://dbpedia.org/resource/Classics	http://aurehal.archives-ouvertes.fr/subject/shs.class
demo	Demography	http://calenda.org/categories.rdf#categorie203 http://dbpedia.org/resource/Demography	http://aurehal.archives-ouvertes.fr/subject/shs.demo
droit	Law	http://calenda.org/categories.rdf#categorie251 http://dbpedia.org/resource/Law	http://aurehal.archives-ouvertes.fr/subject/shs.droit
eco	Economies and finances	http://calenda.org/categories.rdf#categorie236 http://dbpedia.org/resource/Economy	http://aurehal.archives-ouvertes.fr/subject/shs.eco
edu	Education Educación	http://calenda.org/categories.rdf#categorie283 http://dbpedia.org/resource/Education	http://aurehal.archives-ouvertes.fr/subject/shs.edu
envir	Environmental studies	http://dbpedia.org/resource/Environmental_studies	http://aurehal.archives-ouvertes.fr/subject/shs.envir
genre	Gender studies	http://calenda.org/categories.rdf#categorie205 http://dbpedia.org/resource/Gender_studies	http://aurehal.archives-ouvertes.fr/subject/shs.genre
geo	Geography	http://calenda.org/categories.rdf#categorie218 http://dbpedia.org/resource/Geography	http://aurehal.archives-ouvertes.fr/subject/shs.geo
manag	Management	http://calenda.org/categories.rdf#categorie239 http://dbpedia.org/resource/Management	http://aurehal.archives-ouvertes.fr/subject/shs.gestion
hisphilso	History, Philosophy and Sociology of Sciences	http://dbpedia.org/resource/History_of_science http://dbpedia.org/resource/Philosophy_of_science http://dbpedia.org/resource/Sociology_of_scientific_knowledge	http://aurehal.archives-ouvertes.fr/subject/shs.hisphilso
hist	History	http://calenda.org/categories.rdf#categorie228 http://dbpedia.org/resource/History	http://aurehal.archives-ouvertes.fr/subject/shs.hist
info	Communication sciences	http://calenda.org/categories.rdf#categorie271 http://dbpedia.org/resource/Communication_sciences	http://aurehal.archives-ouvertes.fr/subject/shs.info
lang	Linguistics	http://calenda.org/categories.rdf#categorie268 http://dbpedia.org/resource/Linguistics	http://aurehal.archives-ouvertes.fr/subject/shs.langue
litt	Literature	http://calenda.org/categories.rdf#categorie269 http://dbpedia.org/resource/Literature	http://aurehal.archives-ouvertes.fr/subject/shs.litt
museo	Cultural heritage and museology	http://dbpedia.org/resource/Museology	http://aurehal.archives-ouvertes.fr/subject/shs.museo
musiq	Musicology and performing arts	http://dbpedia.org/resource/Musicology http://dbpedia.org/resource/Performing_arts	http://aurehal.archives-ouvertes.fr/subject/shs.musiq
phil	Philosophy	http://calenda.org/categories.rdf#categorie261 http://dbpedia.org/resource/Philosophy	http://aurehal.archives-ouvertes.fr/subject/shs.phil
psy	Psychology	http://calenda.org/categories.rdf#categorie266 http://dbpedia.org/resource/Psychology	http://aurehal.archives-ouvertes.fr/subject/shs.psy
relig	Religions	http://dbpedia.org/resource/Religion	http://aurehal.archives-ouvertes.fr/subject/shs.relig
scipo	Political science	http://calenda.org/categories.rdf#categorie242 http://dbpedia.org/resource/Political_science	http://aurehal.archives-ouvertes.fr/subject/shs.scipo
socio	Sociology	http://calenda.org/categories.rdf#categorie201 http://dbpedia.org/resource/Sociology	http://aurehal.archives-ouvertes.fr/subject/shs.socio
stat	Methods and statistics	http://dbpedia.org/resource/Statistics	http://aurehal.archives-ouvertes.fr/subject/shs.stat

How to use vocabularies to enrich GoTriple ?

- 1 Enriching GoTriple data thanks to vocabularies ?
- 2 The TRIPLE Vocabulary and MORESS
- 3 The annotation and classification GoTriple services based on vocabularies**
- 4 Conclusion & future work
- 5 Questions



Discover
Connect
Collaborate

A decorative graphic on the right side of the slide features three thick, overlapping sticks. One stick is purple, one is red, and one is yellow. They are positioned at an angle, with their ends pointing towards the bottom right corner of the slide.

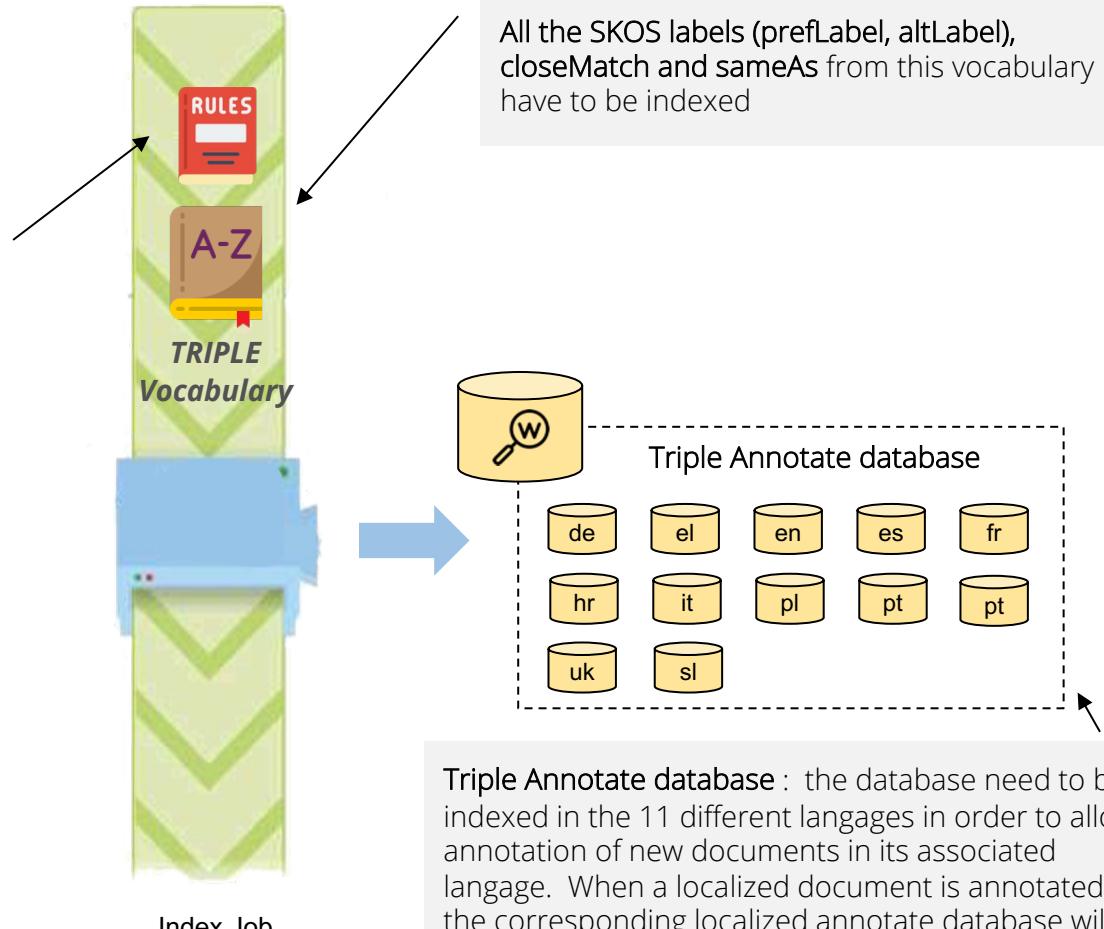
TRIPLE Vocabulary and annotations: build

Annotation

GoTriple uses a **semantic annotation service** that, after being deployed and linked with the **TRIPLE Vocabulary**, tags all TRIPLE keywords in documents in GoTriple.

Semantic rules and applied normalizations:

- Case sensitive or insensitive
Example : "Geography" = "geography"
- Accent sensitive or insensitive
Example : "géographie" = "geographie"
- Inflection sensitive or insensitive
Example : "géographies" = "géographie"
- Stop words removing
Example: "The Geography" => try to find " Geography"
- Inflections
Example: "History Geography" => try to find " Geography History"

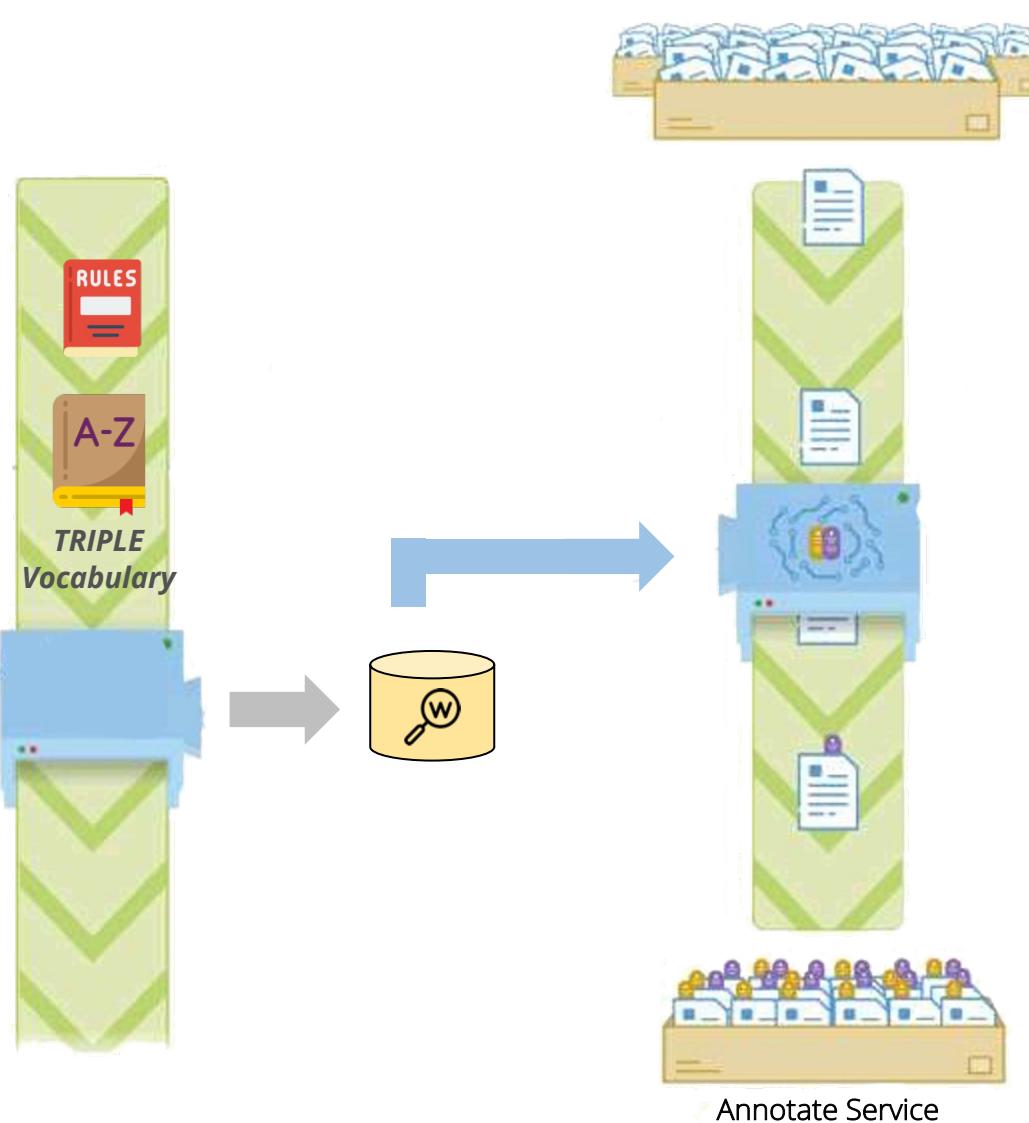


Discover Connect Collaborate



TRIPLE Vocabulary and annotations: serving

Annotation



Documents to annotate

These documents are sent to a **specific API** with 2 inputs :

- Metadata : lang
- Plain text (could be title, description, abstract or full text)

The **Annotate Service** is a script which are able to request the Triple indexed database from a set of metadata and the full text of the document. This script is the link between the annotation database and the API.

The Annotate Service will **analyze** input text to detect if labels (and extended labels based on semantic rules) from the **Triple Vocabulary** are present in the plain text. This detection is based on **exact matches + semantic rules (case/accents/flexions sensitivity)**

Connect
Collaborate

TRIPLE Vocabulary and annotations: challenging points

Annotation

The TRIPLE Vocabulary is built on more than 3,300 concepts. Therefore, the Annotate Service can return many results. We only need to filter the most relevant TRIPLE vocabulary tags between the different matched labels.

Strategy #1

Filter TRIPLE vocabulary tags that are "blacklisted"

Example :

- Concepts with labels that are too generic (Example: "value", "values" ...)
- Concepts with labels built on many stopwords (Example: "the New", ...)
- Concepts with labels based on a combination of the 2 last assumptions (Example: "Other (Philosophy)")

Strategy #2

Filter TRIPLE Vocabulary tags that are not "closed together"

Complete example

Concepts with their top concepts found in a document:

- Concepts: ['Civilization']
- Contracts: ['Law']
- Democracy: ['Civilization']
- Drawing: ['Civilization']
- Fiction: ['Philology']
- Kings and rulers: ['Civilization']
- Plots (Drama, novel, etc.): ['Philology']
- Political science: ['Civilization']
- Politics, Practical: ['Civilization']
- Population: ['Civilization']
- Value: ['Civilization', 'Anthropology']

Filter 1 blacklist:

Blacklisted concepts:

- Value: ['Civilization', 'Anthropology']
- Values: ['Civilization']

Filter 2 "closed together":

Top concepts score:

- ['Civilization: 10', 'Anthropology: 2', 'Law: 1', 'Philology: 2']

We keep group of same-top-level keywords with size ≥ 3 :

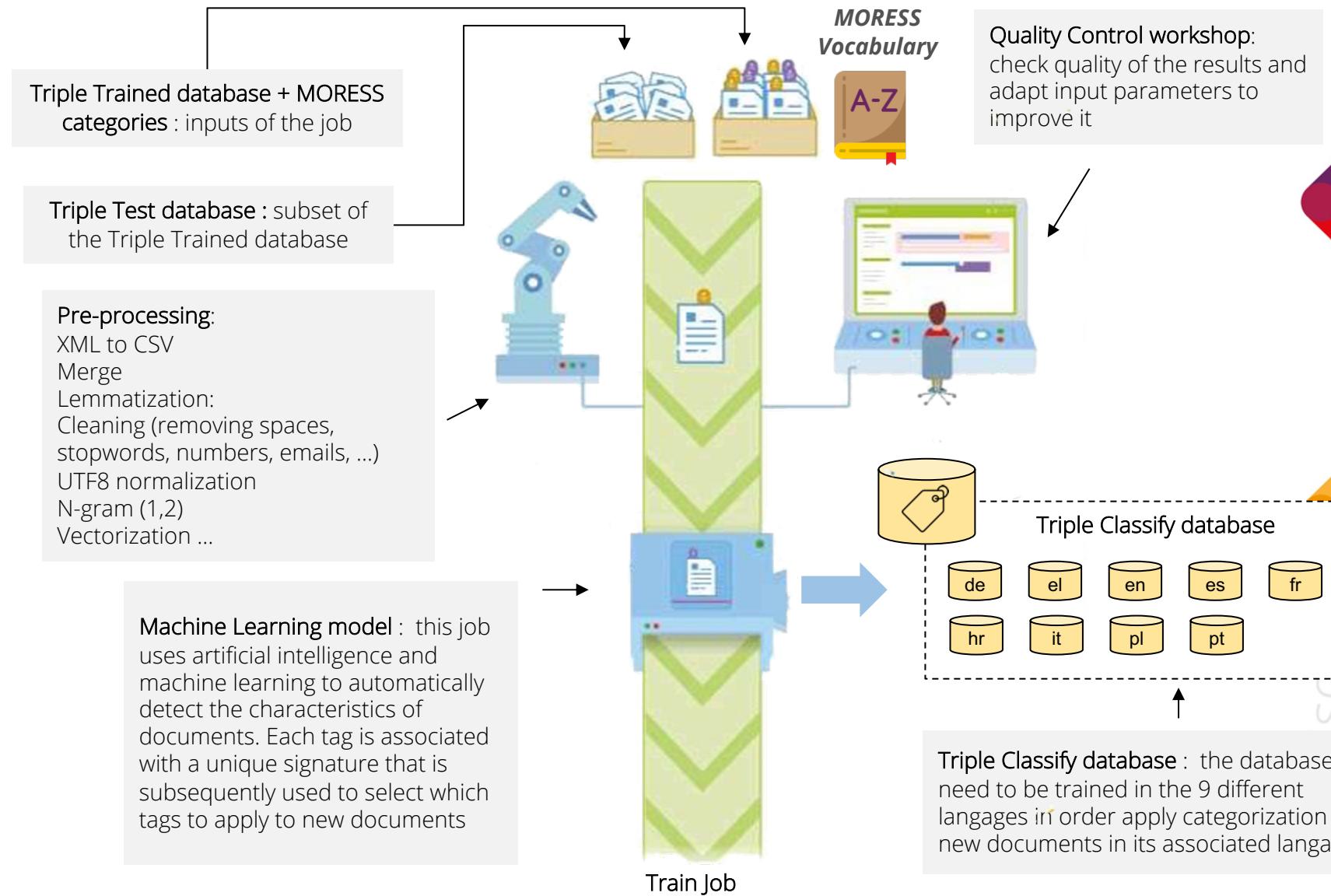
- Final keywords = Biopolitics, Democracy, Political science, "Politics, Practical", Concepts, Drawing, Population



MORESS Vocabulary and classify: build

Classification

GoTriple uses a **semantic classify service** for classification, one or more "discipline", that is the **27 MORESS categories** that have been selected in the TRIPLE project as representatives of the SSH domain

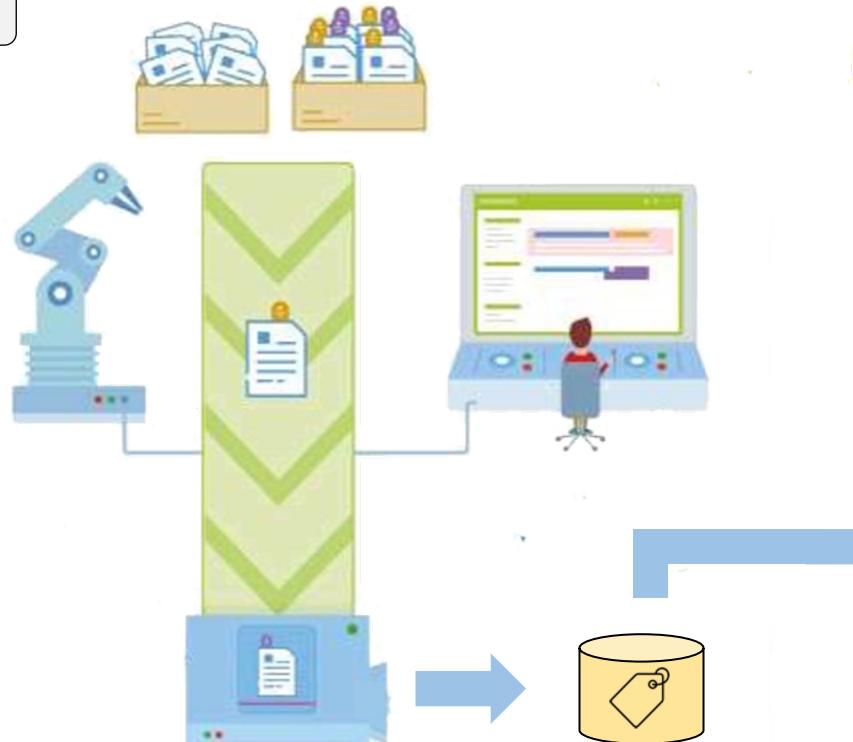


SC
connect
Collaborate
riple

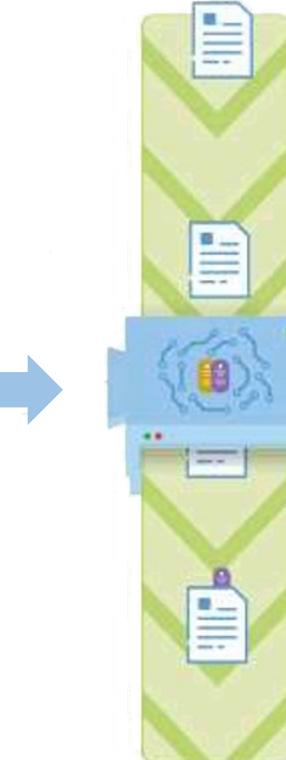
MORESS Vocabulary and classify: serving



Classification



For each document, a score is calculated **for each of the 27 MORESS categories.**



Classify Service

New documents to tag

These documents have to be sent with 3 inputs :

- Metadata : lang
- Theashold : float between 0 and 1
- Plain text (could be title, description, abstract or full text)

The **Classify Service** is a script which are able to request the Triple trained database from a set of metadata and the full text of the document. This script is the link between the prediction model based on MORESS and the API.

The prediction model will analyze its semantic proximity to the different MORESS categories and produce a score.

Discover
Connect
Collaborate

**MORESS
Vocabulary**



Triple

MORESS Vocabulary and classify: challenging points



Classification

1. The Classify Service can return all predictions in every MORESS category. We need to filter only the most relevant categories between the different scores of predictions.

Strategy #1

Filter categories with a score higher than a predefined threshold to maximize the quality rather than the number of categories, even if it means having more documents without categories. The goal is to maximize a user's trust so that they do not lose overall trust for the entire feature.

Strategy #2

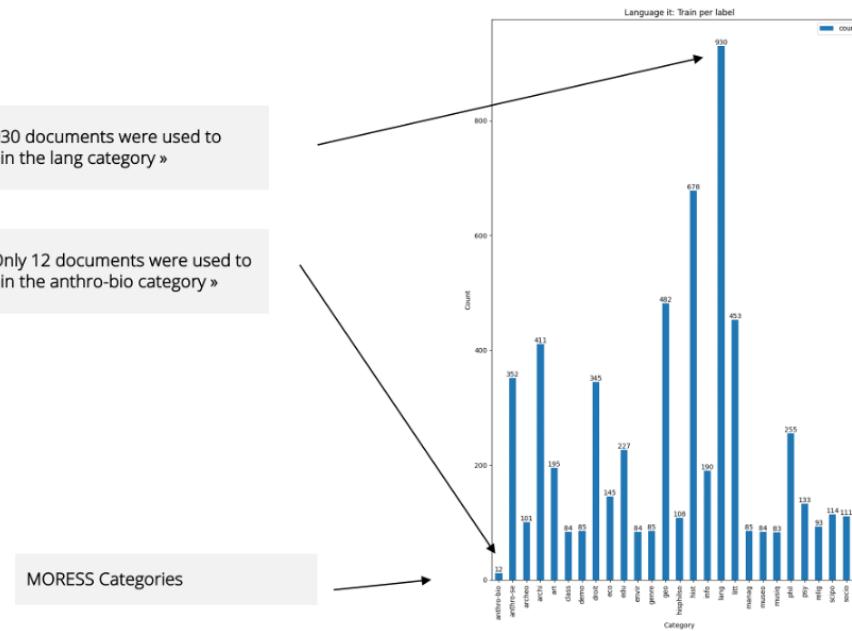
Filter categories that have no more than 2 categories, as this is considered irrelevant in the context of the documentary corpus from SSH.

Some categories are overrepresented or underrepresented in the train model.

Strategy

In the process of evaluation of the train model (the Quality Control Workshop) the MORESS Distribution measure allows us to understand which categories are overrepresented or underrepresented in the train model.

Results with the italiano model



« 930 documents were used to train the lang category »

« Only 12 documents were used to train the anthro-bio category »

MORESS Categories



How to use vocabularies to enrich GoTriple ?

- 1 Enriching GoTriple data thanks to vocabularies ?
- 2 The TRIPLE Vocabulary and MORESS
- 3 The annotation and classification GoTriple services based on vocabularies
- 4 **Conclusion & future work**
- 5 Questions



Discover
Connect
Collaborate

Conclusions and future work

Annotation



- A balance of all input parameters has been found limiting noise in favour of more relevant tags
- The application of the blacklist filter composed of 100 entries allowed the removal of more than 30% of erroneous annotations
- The impact of the application of the "closed together" filter seems to have a positive impact on the overall quality of the results, but remains difficult to assess.

Classification



- Globally, the models (1 model for each language) are efficient: on average, the F1-Score of the 11×27 categories $> 50\%$
- The (relative) perception in terms of user satisfaction is rated "positive" to "very positive". Overall, users have confidence in the discipline displayed to them. This feeling can be partly explained by the use of the confidence threshold and the limit of the number of results to 2.
- The results are unequal between the models (ex: HR has better results than FR)
- The results are unequal between categories (ex: anthro-bio VS lang)

- Use a Machine Learning-based algorithm
- Cross the results of the Annotate Service with the results of the Classify Service

- Review the content of the training base and the training corpus to complete and enrich it (missing fields, representativeness of categories, etc.). This is important work but could considerably improve the MORESS distribution and quality of the results observed.
- Review the confidence threshold and the max limit (=2)

Discover Connect Collaborate



Triple

TRIPLE will be a dedicated service of the OPERAS RI

OPERAS
open access in the european research area through scholarly communication



This project has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement number 863420

CC BY 4.0 International Licence

Julien Homo – julien.homo@foxcub.fr

Collaborate

ple