Corpus Viewer platform relies on Natural Language Processing (NLP), Machine Learning (ML) and Machine Translation (MT) to analyze structured metadata and unstructured textual data in large document corpora. The platform allows the decision maker and the policy implementer the possibility of analyze R&D&i information space (mainly patents, scientific publications and public aids) for evidence and knowledge-based policy making and implementation. It relies, among other techniques, on topic modeling and graph analysis.

The development of Corpus Viewer started in 2016 and is still progressing based on the collaboration of several subcontracted University research groups and companies. Corpus Viewer on its 1.0 version, is currently used by three public administrations: SEAD (Ministry of Economy), the Spanish Foundation for Science and Technology (FECYT) and the State Secretary for University and Research, Development and Innovation (SEUIDI) at the Spanish Ministry of Science.

Although Corpus Viewer is a generic platform that can be exploited with virtually any collection of text documents, the current deployment of the platform mainly hosts R&D related text corpora:

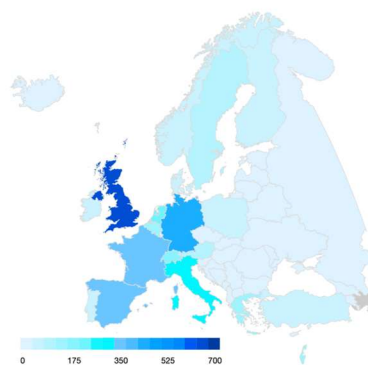| Corpus | Docs in corpus | Time Frame |
|---|---|---|
| Projects funded by the Spanish State Plan for R&D | 110 K | 2004 - 2016 |
| European R&D projects (CORDIS) | 78 K | 1984 - 2018 |
| American R&D projects (NSF) | 150 K | 1985 - 2017 |
| American Health R&D projects (NIH) | 1.8 M | 1983 - 2017 |
| Patent applications (PATSTAT) | 90 M | 1898 - 2017 |
| Scientific papers with contributions from authors affiliated to a Spanish institution (SCOPUS) | 680 K | 2006 - 2018 |

These data sources are processed to assist in the definition and implementation of R&D&i public policies through a set of functionalities allowing to:

1. compare R&D&i funding and knowledge areas in different geographic regions,
2. identify competitive advantages between countries, regions, organizations,
3. identify R&D&i knowledge areas, as well as their emergence, evolution and even hybridization with other knowledge areas (it provides also metadata aggregation and BI type dashboard visualization),
4. R&D agent (organization, researcher and firm) profiling and,
5. assist in the assessment of the impact of public policies by tracking the outputs of grants, short and long term outcomes in terms of lead-lag.
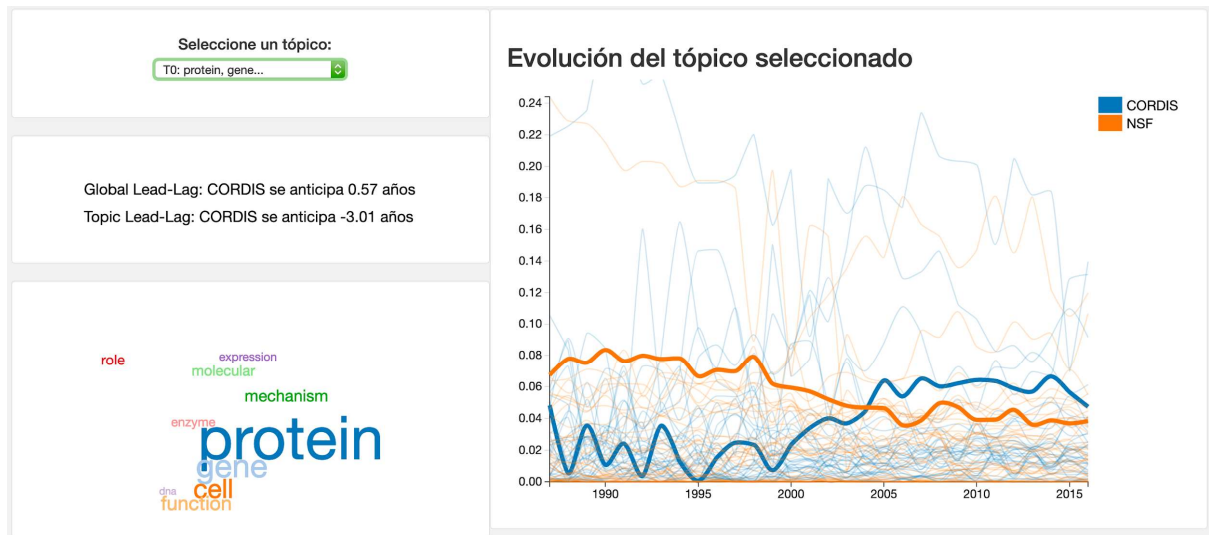
Corpus Viewer also provides tools for policy implementation, in particular for the selection of evaluators or the retrieval of relevant documents (patents, scientific publications, R&D aid grants and proposals, etc) for innovation evaluation. Furthermore, it is used for plagiarism detection, identification of cases of double funding and fraud in aid grants and proposals submitted for national funding.
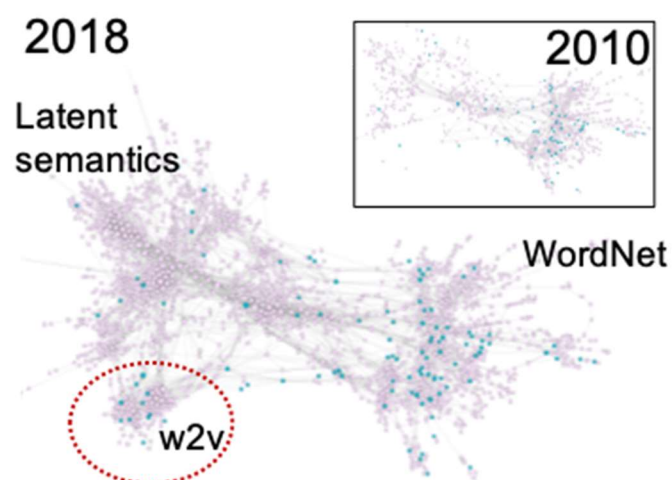
## Corpus Viewer functionalities

1. Scalable NLP pipeline and Machine Translation of large volumes of documents.
2. Automatic classification of documents according to available taxonomies using deep learning networks.
3. Topic modeling analysis of document collections and topic inference for new documents.
4. Information retrieval system based document similarity.
5. Optimized columnar and textual indexing for efficient searches and query metadata filtering.
6. Tracking of semantically alike documents (semantic clusters) (emergence, evolution and hybridization with other clusters.
7. Dynamic topic analysis and temporal thematic evolution. Temporal analysis by areas of knowledge, lead-lag between different types of document corpus.
8. Topic-enhanced dashboards to analyze R&D&i distribution by metadata, including geographic area.
9. Automatic profiling and disambiguation of R&D&i key players based on their R&D production.
10. Analysis of collaboration networks

## Corpus Viewer Facilitating Techniques

Unlike traditional statistical description methods, our approach does not require predefined taxonomies or controlled vocabularies. Analysis techniques are based on underlying topics, conceptual indexing of documents, word Embeddings and other techniques of documentary representation.

1. Scalable tokenization, PoS tagging, lemmatization, disambiguation and wikification for English and Spanish languages.
2. Automatic translation (ES-EN), NMT based.
3. Topic modeling, including: static models (Latent Dirichlet Allocation, LDA, CTM), dynamic topic models, hierarchical LDA, and recursive LDA
4. Textual search and document similarity implemented for topic, bag of words, and word Embeddings document representations
5. Analysis of graphs; modularity, distances between clusters and centrality calculation

# Corpus Viewer Use Cases

Potential users of current Corpus Viewer platform include:

1. R&D policy makers
2. Managers and coordinators of R&D programs (policy implementation)
3. Grant evaluators
4. R&D agents: researchers, Public Research Institutions, Companies

| CASO DE USO | Perfil de usuario | | | | |
|---|---|---|---|---|---|
| | Decisor | Gestor Ayudas | Evaluador | Investigador | Empresa |
| DISEÑO POLÍTICAS PÚBLICAS (seguimiento, prospectiva, planificación) | X | X | X | X | X |
| Herramientas de soporte a EVALUACIÓN (evaluadores, similitud de doc., clasificació taxonomías, estimación innovación) | | X | X | X | X |
| SISTEMA DE ALARMAS (plagio, patrones de fraude) | | X | X | | X |
| SISTEMA DE RECOMENDACIÓN (evaluadores, licitación, formación, cruce, conocimiento implícito) | X | X | | X | X |

# Corpus Viewer Demonstrators

Corpus Viewer access is restricted to authorized users. However, within the project several demonstrators [Enlace a la Página web de demostradores de IntelComp] have been developed that provide access to part of the functionality implemented in the platform.