# Making Sense of Massive Amounts of Scientific Publications:

# The Scientific Knowledge Miner Project

Francesco Ronzano, **Ana Freire**, Diego Saez-Trumper, Horacio Saggion

**upf.** **Universitat Pompeu Fabra** *Barcelona*

# 20 seconds... 1 paper

The Scientific Knowledge Miner Project

# Information Overload (scientific repositories)

# Information Overload (scientific repositories)

**24,6M** PubMed

**90M** WEB OF KNOWLEDGE℠ THOMSON REUTERS

**1M** arXiv.org

**57M** Scopus ELSEVIER

Sometimes between 2017 and 2021, more than half of the papers available globally are expected to be published as Open Access articles.

Lewis, David W. **"The inevitability of open access**."
College & Research Libraries 73.5 (2012): 493-506.

# The peculiarities of research publications

**TITLE**

**ABSTRACT**

**(SUB)SECTION**

**CAPTION**

**BIBLIOGRAPHIC ENTRY**

# Scientific publications: claims

In order to take full advantage of the knowledge present in scientific publications proper **semantic indexing**, **search** and **content aggregation** approaches, are required.

Benefits:

- Search of new information on specific scientific problems
- Semi-automatic assessment of papers and research proposals
- Hypothesis formulation
- Tracking of scientific and technological advances
- Scientific intelligence
- Assisted report and review writing
- Question answering
- …

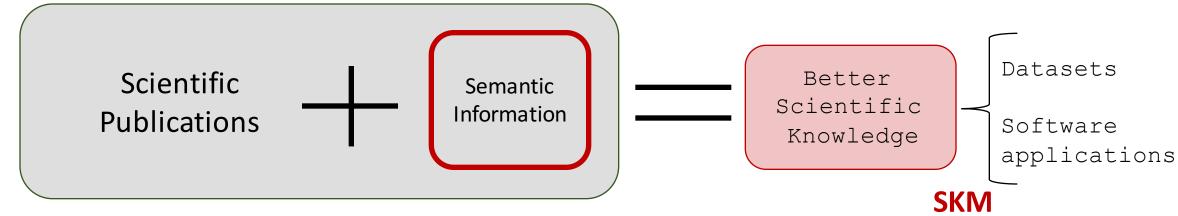# The Scientific Knowledge Miner Project (SKM)

Facilitate the extraction of knowledge from scientific publications across many disciplines.

Improve a variety of use cases such as:
- Citation Characterization
- Citation Recommendation
- Summarization
- …

➤ KEY: Papers are enriched with **structural, linguistic** and **semantic information**
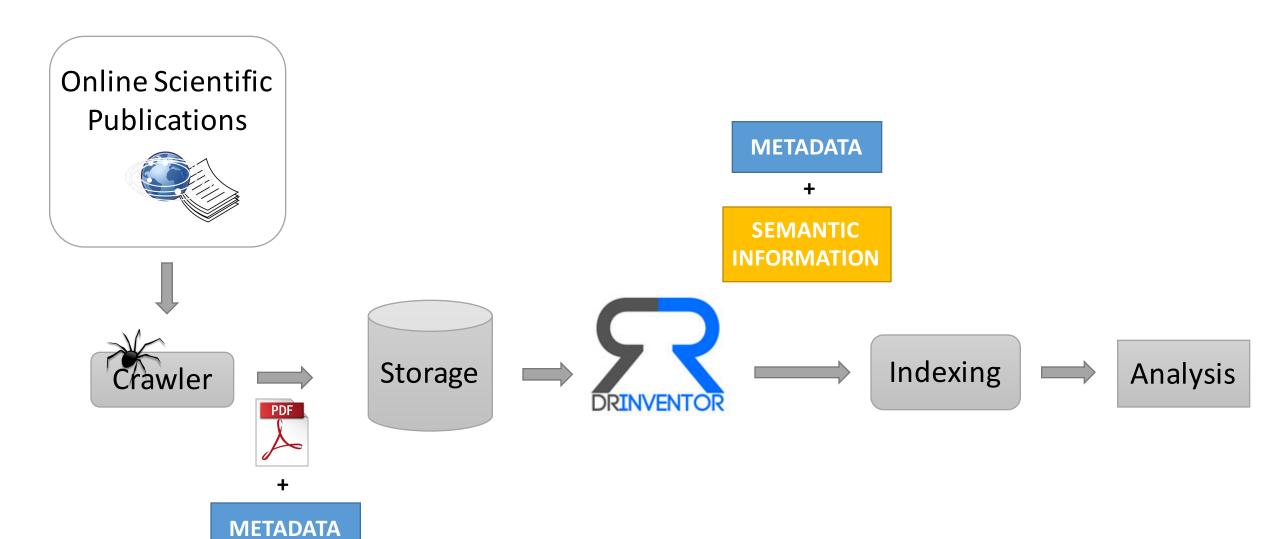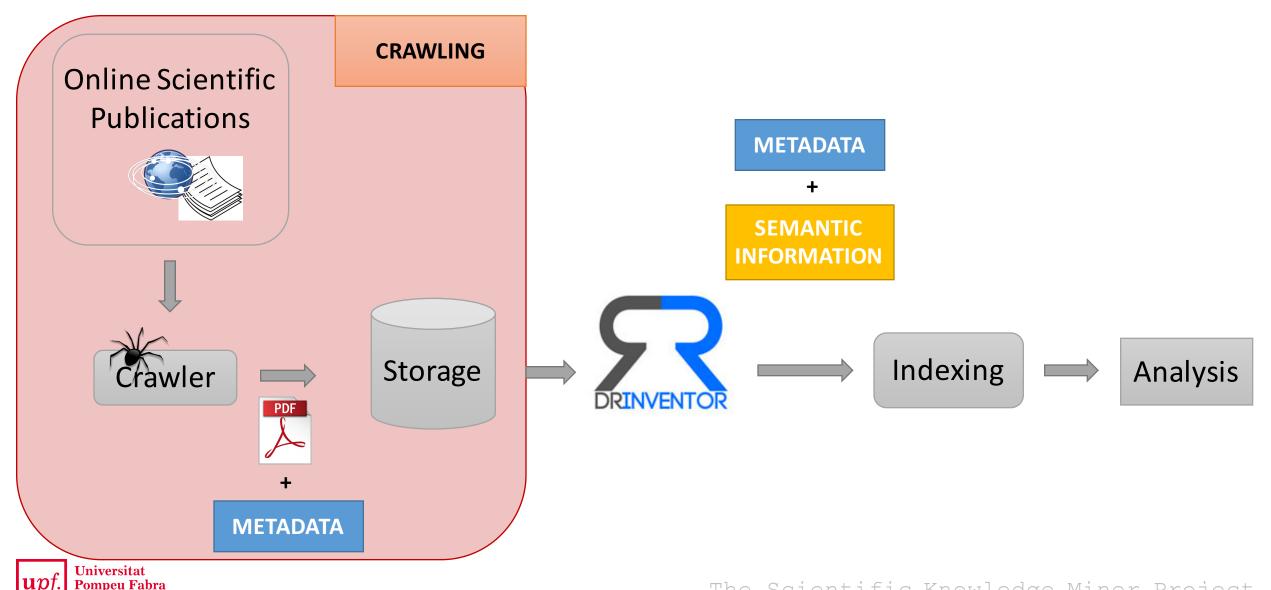
# The Scientific Knowledge Miner Project (SKM)

The SKM approach to the analysis of scientific literature:

- Relies on a finer-grained analysis of the contents of publications

- Is grounded on the automated characterization of a varied set of semantic aspects of papers, including the rhetorical structure or the purpose of citations.

# The Scientific Knowledge Miner Project (SKM)

# The Scientific Knowledge Miner Project (SKM)

# Crawling

# The Scientific Knowledge Miner Project (SKM)

# The Scientific Knowledge Miner Project (SKM)

# Dr. Inventor Text Mining Framework

- Integrate and customize **text mining tools** and **on-line services** to enable and ease a wide range of scientific publication analyses

- Papers are enriched with **structural**, **linguistic** and **semantic information**

**http://backingdata.org/dri/library/**

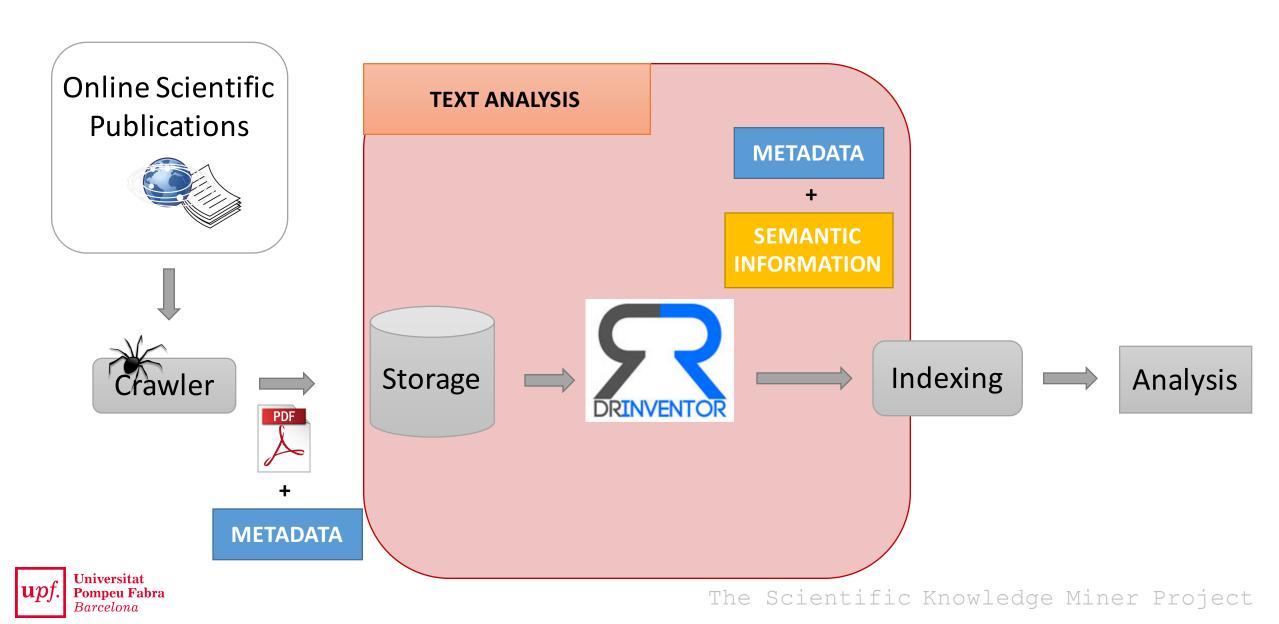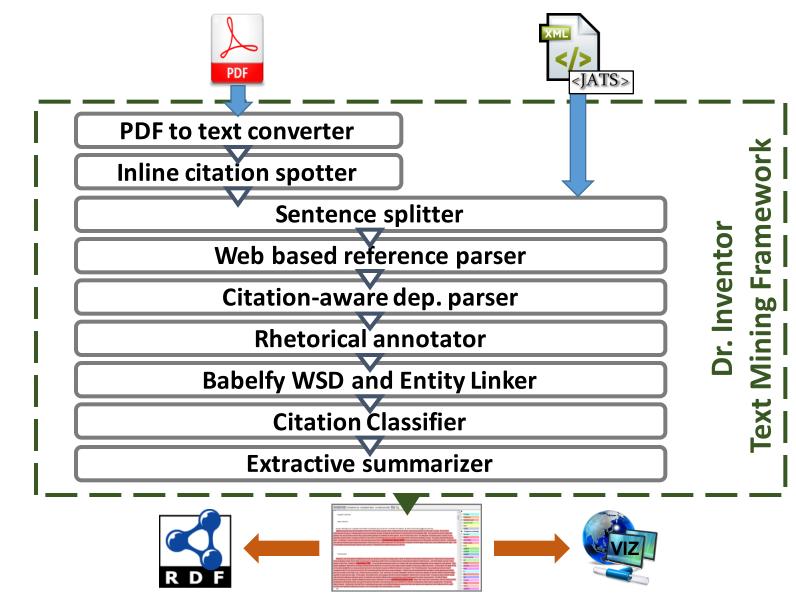- Self-contained **Java 8** library managed by **Maven™**

- Focused on **textual content**

- Relying on a **shared data model** (java classes) to represent a paper

- Exposing a **convenient API** to access the mined information
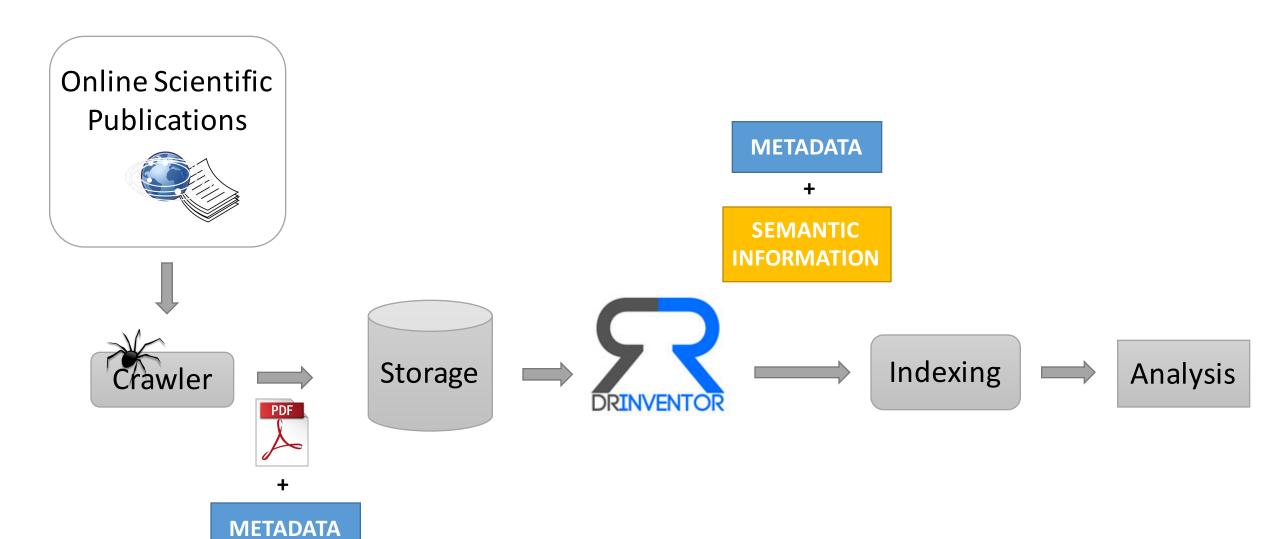
- Based on **GATE** *general architecture for text engineering* to manage **textual annotations**

# Dr. Inventor Text Mining Framework

# The Scientific Knowledge Miner Project (SKM)

# The Scientific Knowledge Miner Project (SKM)

# Indexing

# The Scientific Knowledge Miner Project (SKM)

# The Scientific Knowledge Miner Project (SKM)

# Analysis

# Use Case 1: Citation Characterization

Experiment new metrics: what do others say about one paper?



Enrich citation counts with semantics

**CITATION PURPOSE**

| Criticism |
| :---: |
| **Comparison** |
| **Use** |
| **Substantiation** |
| **Basis** |
| **Neutral** |

**+ 17 sub-purposes**

# Use Case 2: Citation Recommendation

Recommend similar papers / authors

**SENTENCE RHETORICAL CATEGORY**

Background

Approach

Challenge

Outcome

Future Work

**+ 3 sub-categories**

Some alternative phrase alignment approaches have been developed, which do not rely on the Viterbi word alignment. Both (Marcu, 2002) and (Zhang, 2003) consider a sentence pair as different realizations of a sequence of concepts. These alignment approaches segment the sentences into a sequence of phrases.

CHALLENGE
BACKGROUND
HYPOTHESIS

# Use Case 3: Scientific Document Summarization

Extractive summarization

Some alternative phrase alignment approaches have been developed, which do not rely on the Viterbi word alignment. Both (Marcu, 2002) and (Zhang, 2003) consider a sentence pair as different realizations of a sequence of concepts. These alignment approaches segment the sentences into a sequence of phrases.

Summary:
Some alternative phrase alignment approaches have been developed, which do not rely on the Viterbi word alignment.
These alignment approaches segment the sentences into a sequence of phrases.

**SENTENCE SUMMARY RELEVANCE (1 to 5 ratings)**

and

**HAND-WRITTEN SUMMARY**

The Scientific Knowledge Miner Project

# Conclusions and future work

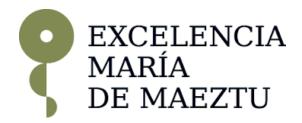Scientific Knowledge Miner (SKM) aims at facilitating the extraction, aggregation and navigation of knowledge from scientific publications.

- Consolidate the SKM publication mining infrastructure
- Exploit the semantics of papers to perform large scale investigations of:
  - Alternative metrics to evaluate a paper based on citation semantics
  - Semantically motivated recommendation of scientific publications
  - Summarization of scientific literature

# Acknowledgements

# Making Sense of Massive Amounts of Scientific Publications:

# The Scientific Knowledge Miner Project

{francesco.ronzano, ana.freire, diego.saez, horacio.saggion}@upf.edu

**Universitat Pompeu Fabra**
*Barcelona*