openM1N7ED

Sophie Aubin - INRA

# **OpenMinTeD**

## Text mining services for e-infrastructures

CAPSELLA Open Data Workshop, Chania, 2 June 2017

# Amounts of scientific texts...

**1 paper/sec**

**90%** of papers **never** cited*

**120,000** papers published on a single taxon **"Zea mays"**

**Publications**

**50%** of papers **never** read by anyone than its authors, referees, journal editors*

+ reports, patents, books, surveys, news, etc.

openM1N7ED

CAPSELLA Open Data Workshop, Chania, 2 June 2017
This is where the footer goes

# TDM services

**Indexing** documents and datasets

Entities **identification** and **normalisation** against reference data

Information **extraction**: from semi-structured to structured

**Semantic/lexical** resource **acquisition** from texts

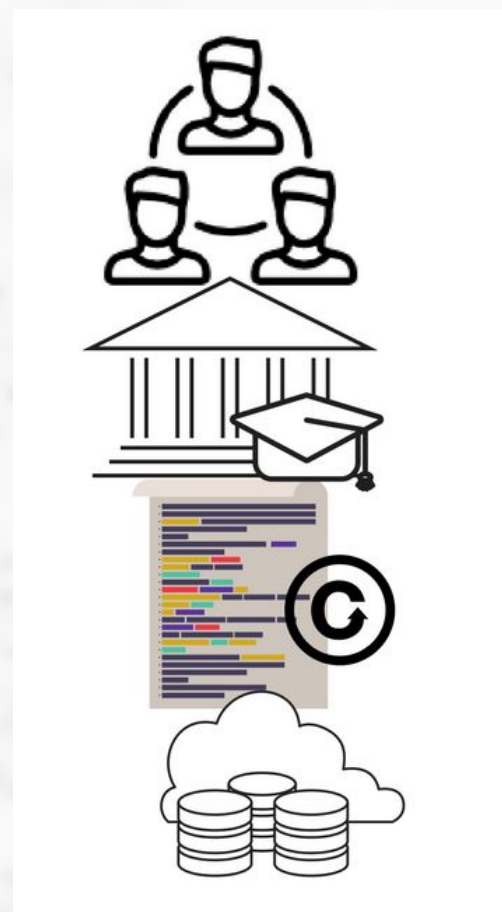# However, …

## Multitude of solutions catering for different

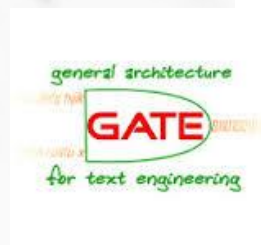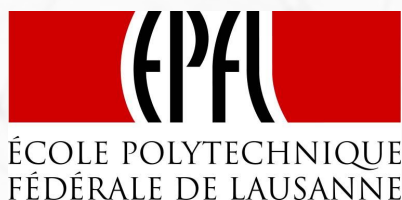| Text Types | Domains | Tasks | Languages |
|---|---|---|---|
| Newswire | Agriculture | Translation | English |
| Scientific Literature | Health | Knowledge acquisition | French |
| Tweets/blogs | Biology | Semantic Search | German |
| Patents | Social Sciences | Question Answering | Spanish |
| Clinical/medical records | Environment | Sentiment Analysis | Portuguese |
| Textbooks, monographs | …. | Summarization | Italian |
| Online forums | | Knowledge Discovery | Polish |
| …. | | …. | …. |

## Creating a fragmented landscape

# Then comes openM1N7ED

**Duration**: 3 years (2015-2018), **16 Partners**

- research groups in text-mining
- content providers
- a data center
- a library association
- legal experts
- community related partners
- SMEs



openM1N7ED
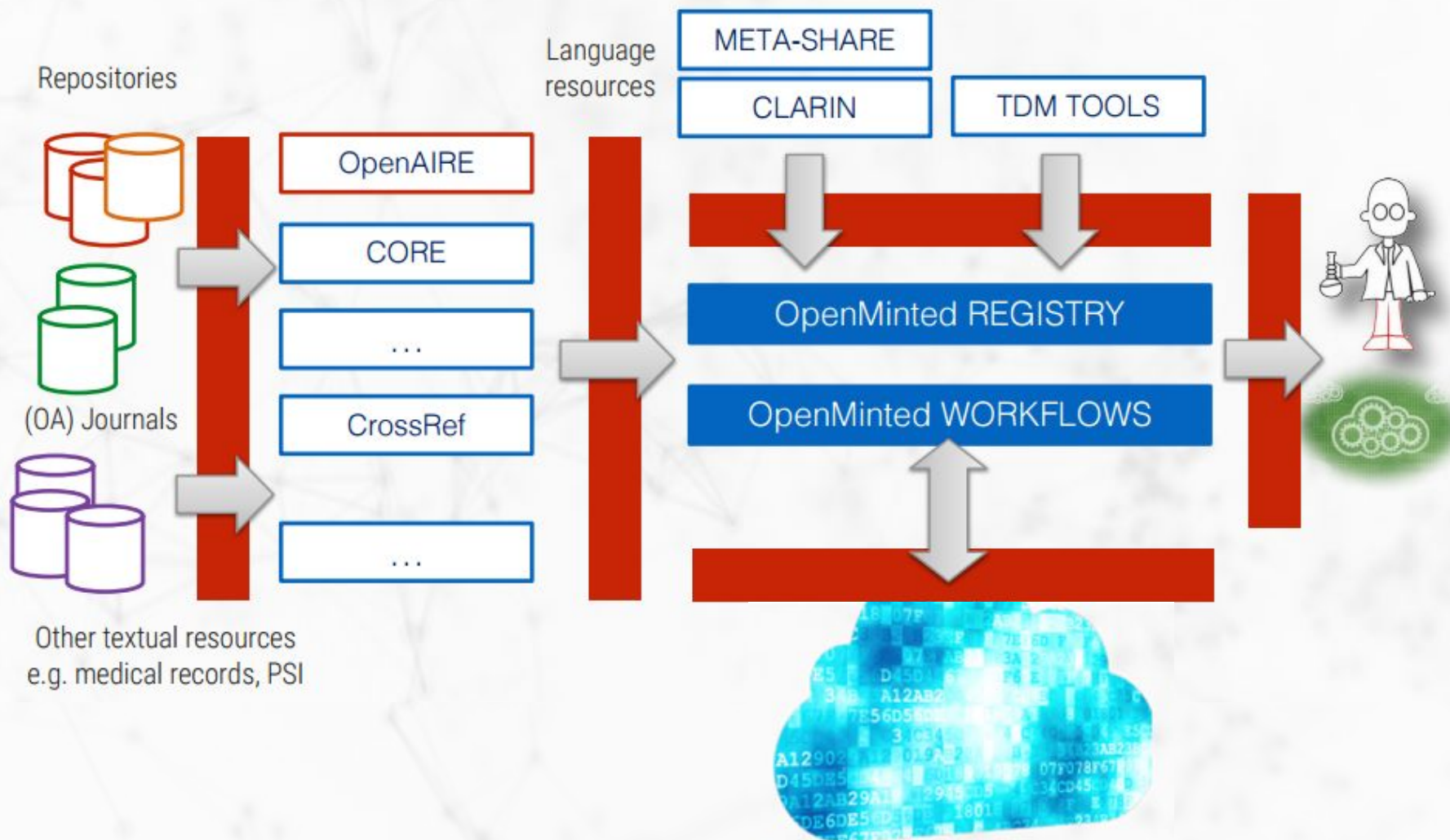
# Partners of openM1N7ED

# What is openM1N7ED?

openM1N7ED is a **platform** that works as an **infrastructural service** of the wider **research ecosystem**

openM1N7ED

European Commission

# Our Services

**1** Discover TDM Services and tools

**2** Feed with home content or easily get texts from (OA) content hubs

**3** Pick adequate knowledge resources

**4** Build your own service/applications – Combine components into a Workflow

**5** Share and Re-Use

CAPSELLA Open Data Workshop, Chania, 2 June 2017

openM1N7ED

# How does all this bind together?

# OpenMinTeD Scientific/Business Applications

**Data & content repositories/registries**
For Data analysts, Scientists
**Curation tools**
For Data providers
**Analytical tools**
For Scientists-policy makers
**Decision tools**
For Farmers, SMEs
**Knowledge acquisition tools**
For Scientists, Ontologists

openM1N7ED

CAPSELLA Open Data Workshop, Chania, 2 June 2017
This is where the footer goes

European Commission

# Example: WheatIS / gnpIS

**Application:** federated search of genomic and genetic data for wheat

**Added value**: direct access from data to related scientific articles
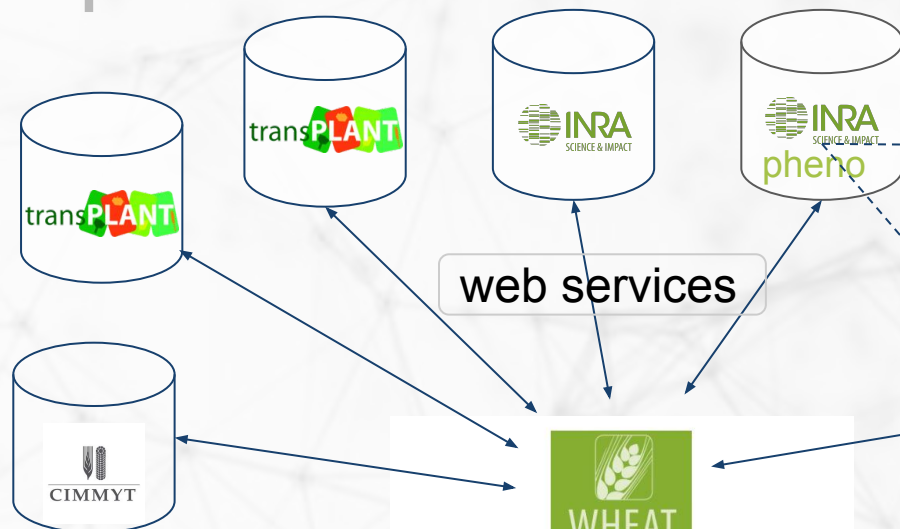
**Objects**: taxa, genes, markers, phenotypes and varieties

**Challenges**: naming heterogeneity & scale variety between textual and experimental data
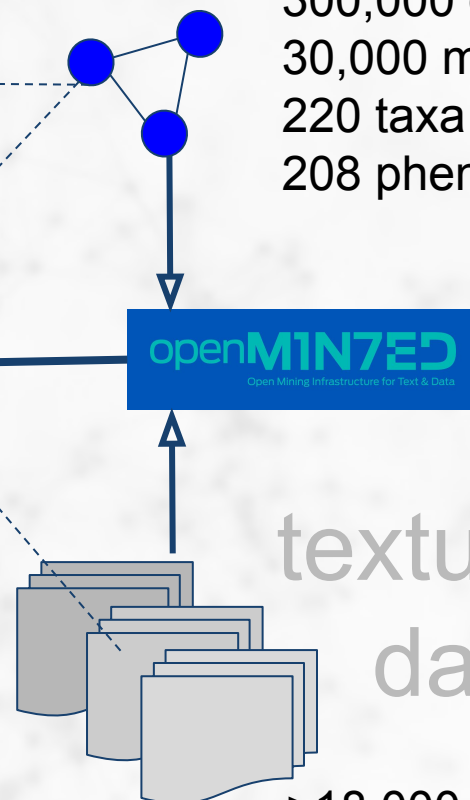
# data - text interoperability



experimental data

Knowledge

300,000 genes
30,000 markers
220 taxa
208 phenotypes

semantic interop.

web services

textual data

>18,000 articles

# Example: living conditions of food-related micro-organisms

**Application:** study and characterize the microbial biodiversity of food ecosystems (dairy or meat products, fish, wine, etc.)

- risk management
- food quality improvement



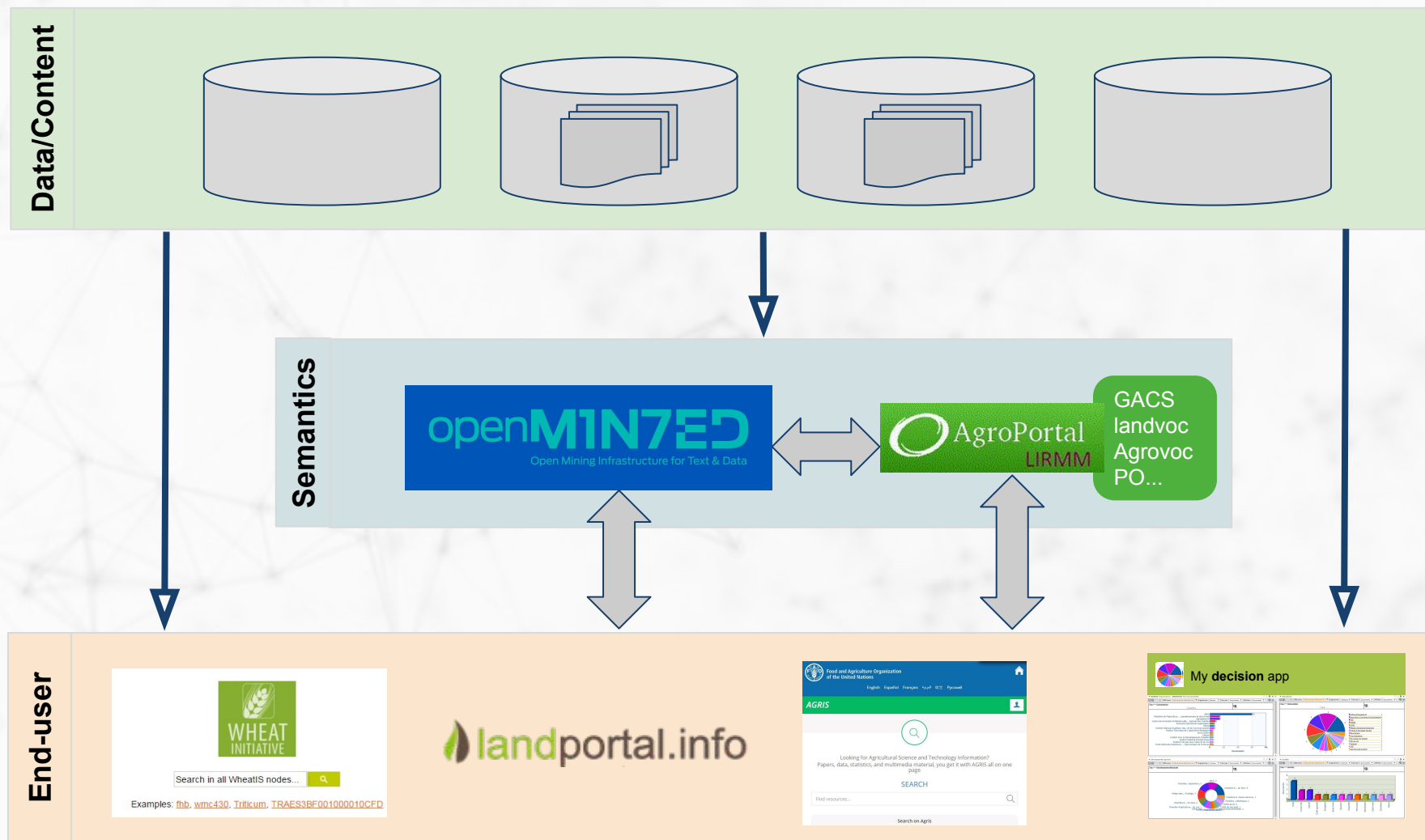**Added value**: completion of knowledge databases with info from the literature



**Objects**: bacteria, habitats and phenotypes

**Challenges**: heterogeneous data integration, object identity

openM1N7ED

CAPSELLA Open Data Workshop, Chania, 2 June 2017

# OpenMinTeD for agriculture



CAPSELLA Open Data Workshop, Chania, 2 June 2017

# Text mining

- creates value from texts
- creates even more value if the results are linked to other data

This relies on shared semantics, standards & protocols

This requires less tech competencies & resources thanks to the common e-infra openM1N7ED

openM1N7ED

# openM1N7ED

# THANK YOU!

Sophie Aubin
sophie.aubin@inra.fr

twitter.com/openminted_eu

facebook.com/openminted

bit.do/openmintedlinkedin

vimeo.com/openminted

bit.do/openmintedplus