# What can semantic text-mining do for food quality improvement?

*Claire Nédellec, Robert Bossy*

*TDM: Unlocking a goldmine of information*
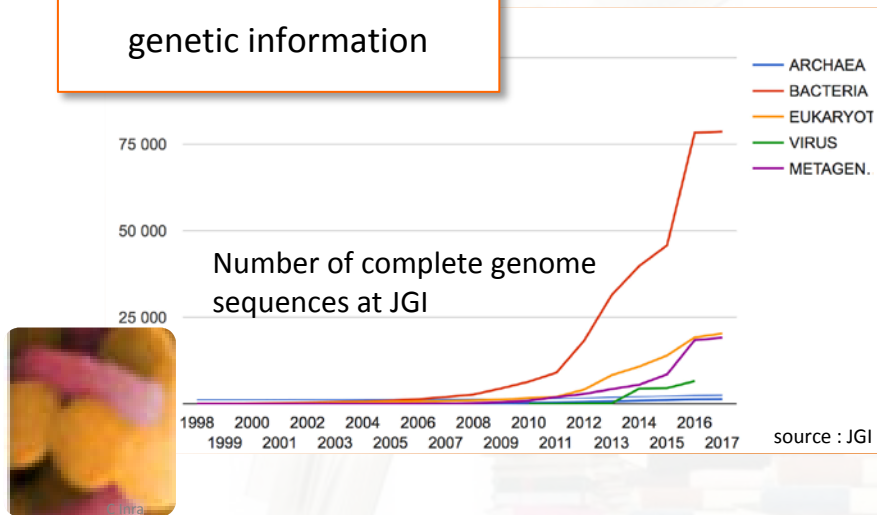**Open Science Fair, Sept. 2017 Athens**

# Microorganisms, food and scientific literature

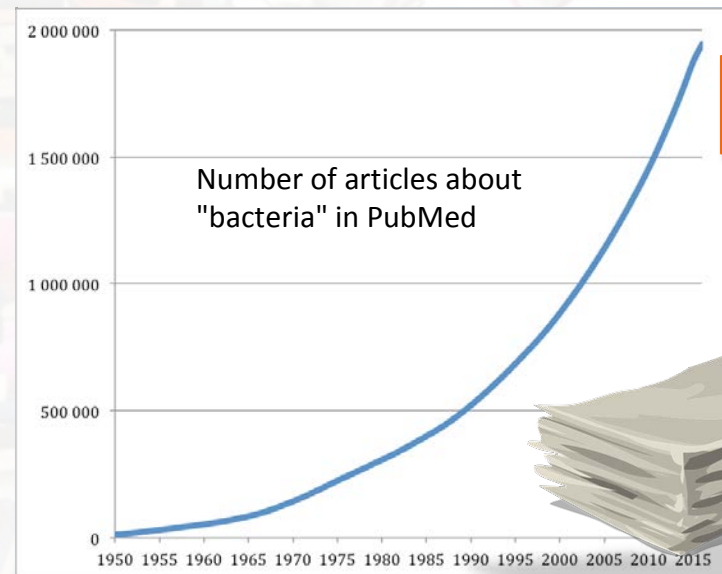Billions of microorganims everywhere.

A critical role in all aspects of our life, *e.g.* food processing.

Researchers study their ecosystems and their genetics for better understanding, control, and use.
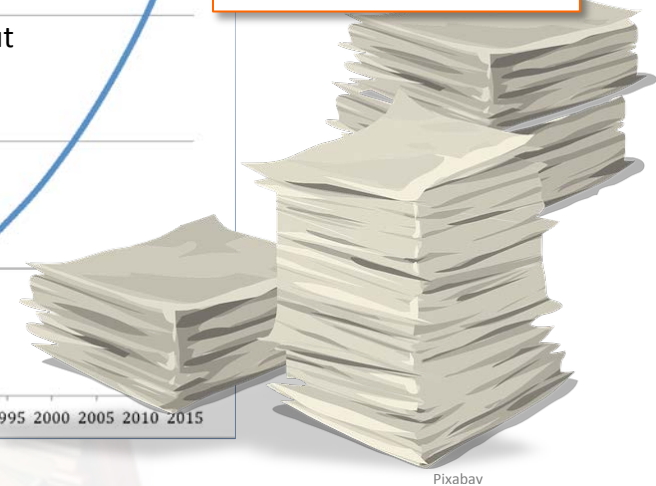
Exponential growing of genetic information

Number of complete genome sequences at JGI

- ARCHAEA
- BACTERIA
- EUKARYOT
- VIRUS
- METAGEN.

75 000
50 000
25 000

1998 1999 2000 2001 2002 2003 2004 2005 2006 2007 2008 2009 2010 2011 2012 2013 2014 2015 2016 2017

source : JGI

Ecosystem, habitats, properties described in millions of free text

Number of articles about "bacteria" in PubMed

2 000 000
1 500 000
1 000 000
500 000
0

1950 1955 1960 1965 1970 1975 1980 1985 1990 1995 2000 2005 2010 2015
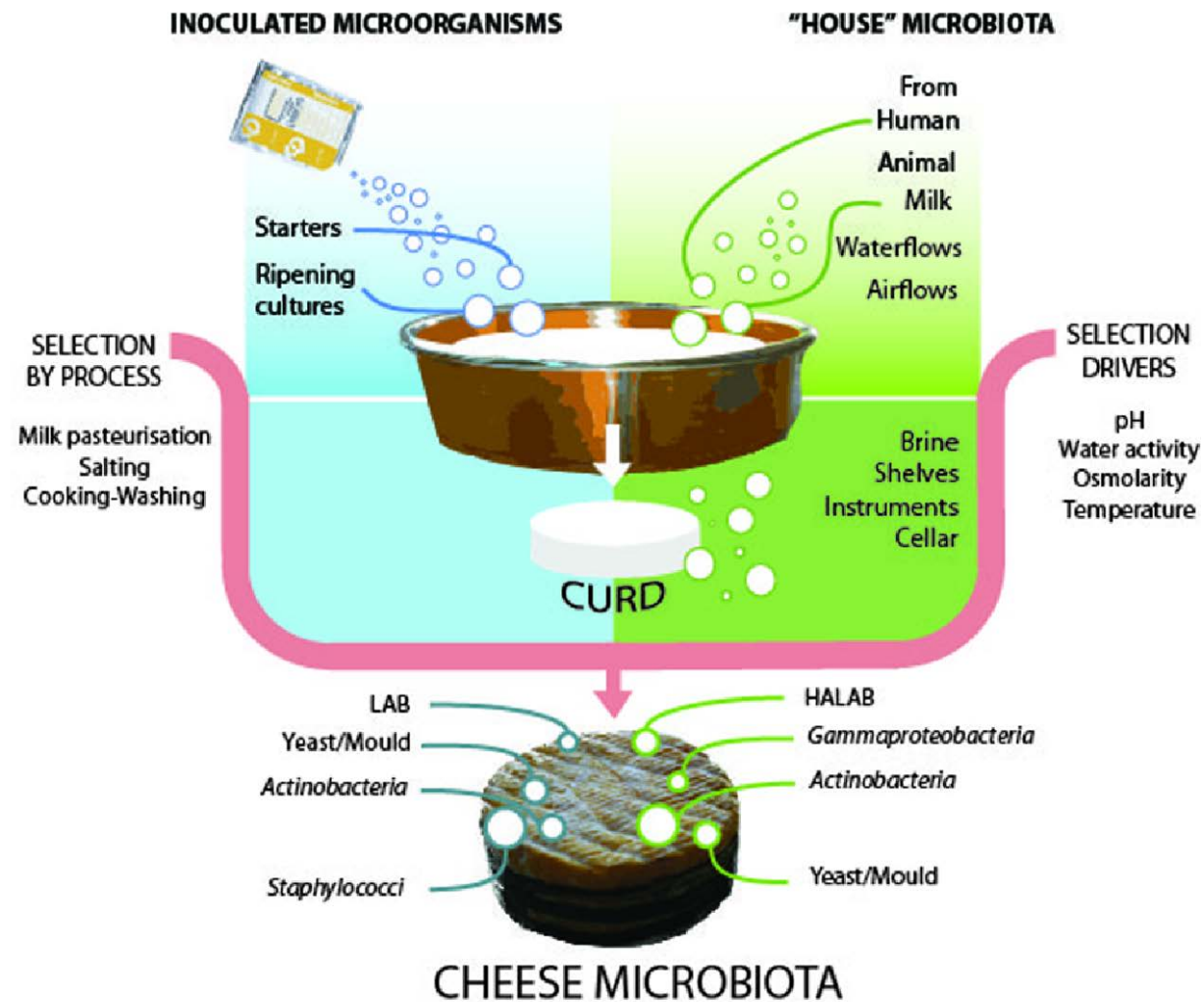
... and of publications

Pixabav

*OntoBiotope*, a semantic text-mining service

A shining example in cheese processing

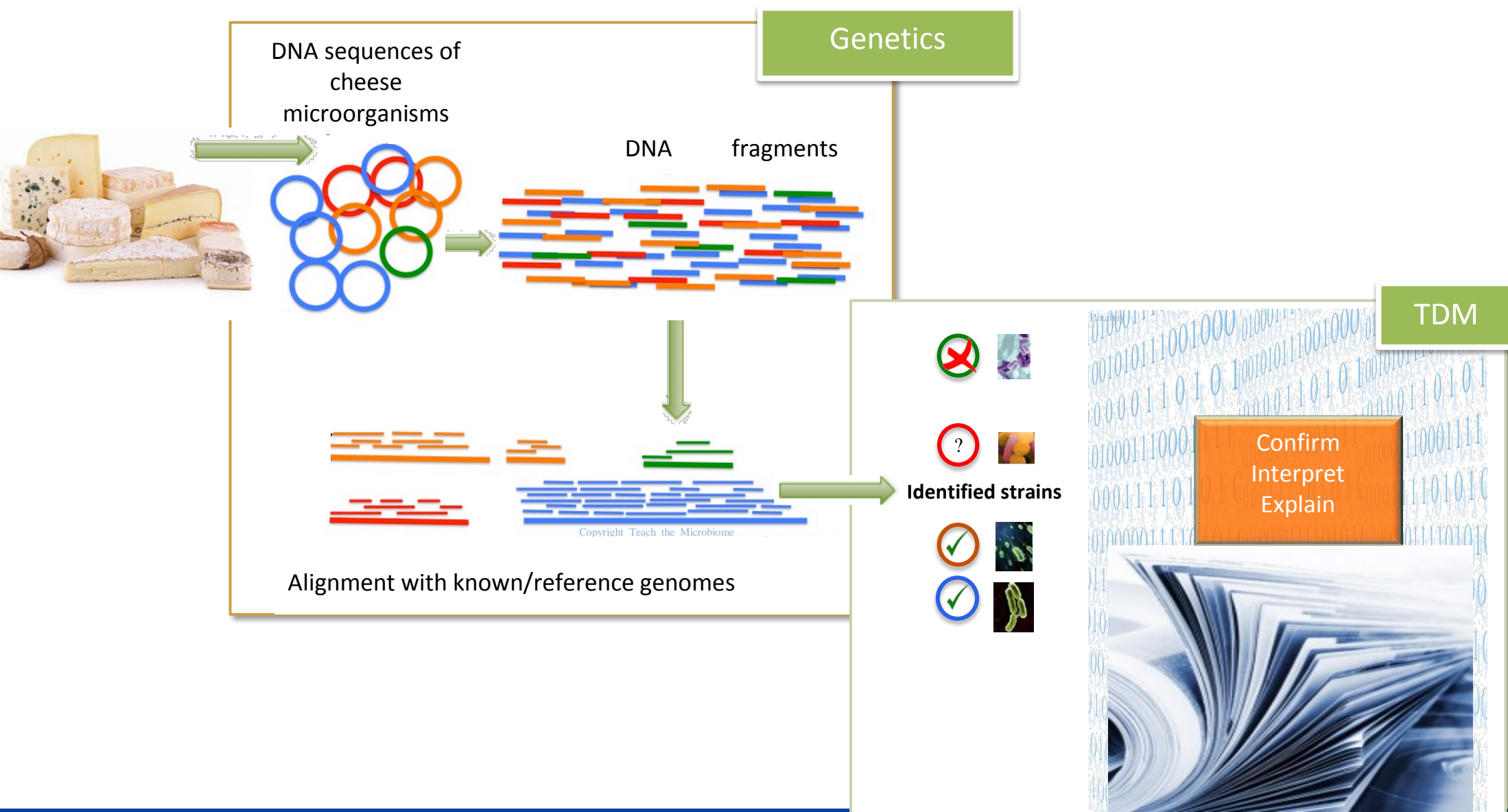An European Open Science perspective

Opportunities and barriers

# Microorganisms in my cheese?



Irlinger et al., FEMS Microbiol Lett (2015) 362 (2).

# DNA identification of microorganisms

DNA sequences of cheese microorganisms

DNA       fragments

Genetics

Alignment with known/reference genomes

Copyright Teach the Microbiome

Identified strains

TDM

Confirm
Interpret
Explain

Metagenomics analysis of hundreds of French and Italian cheese samples

Identify microorganisms to understand and control their presence,
Improve quality of food product and design new ones

Inra - Cniel project
*FoodMicrobiome Transfert*

cheese samples

**among 400 strains**

**300** very frequent and well-known strains

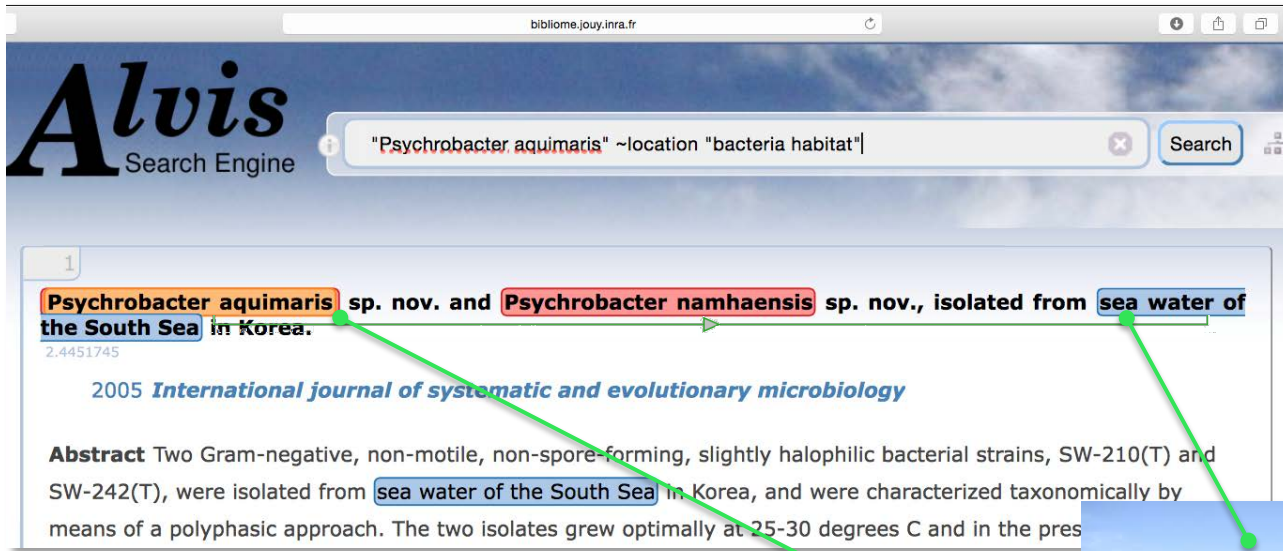**100** strains, little known and varyingly present

strains

Copyright © INRA

**Psychrobacter aquimaris**
ER15_174_BHI7

gorgonzola, roquefort, époisse, toscanello, st nectaire (very frequent), tomme, bleu

# TDM explanation of Psychrobacter presence



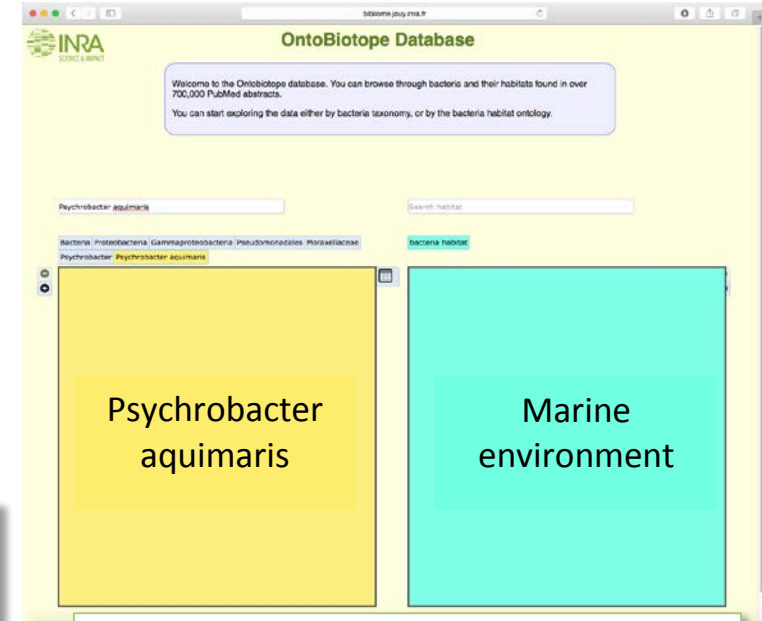TDM OntoBiotope analyzed all PubMed information
- 1,16 millions de documents
- 3,63 millions of relations bacteria-habitat assigned to 2000 categories !
- prediction performances

Psychrobacter aquimaris

Marine environment

*Psychrobacter aquimaris* comes from sea environment

The researcher understands :
the added salt brings the bacteria in the cheese

# Many other good reasons to study food microbiome

**Industrial interest**

Better understanding and control of microbiome role in food process

Food innovation : transformation, preservation, antibiotic or nutrient production, flavour/taste

**Public health**

Better control of food spoilage and safety by explaining the microbial source and adaptation.

Improve intestinal microbiota by food intake.

**Fundamental research**

Better understanding of microorganism life and adaptation

Ecosystem dynamics

In food and related environments

# A TDM based-service for food microbiome study

**The TDM-based service** developed by INRA

- automatically extracts information
- from massive amount of documents in all microbiology domains
- interoperable with experimental data using shared ontologies.

**What information**

- *who* is living: taxa, strain, species, families
- *where*: habitats of all kinds
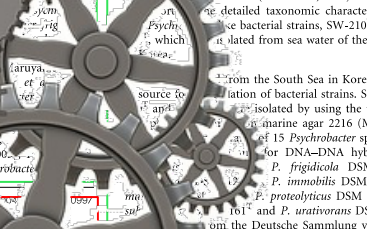- *how*: environment requirements, phenotype

Q?: bacteria livesin food and human

**Abstract** Enterobacter sakazakii is an emerging foodborne pathogen associated with meningitis, necrotizing enterocolitis, and sepsis in infants. One of the main transmission vehicles is the

**Probiotic** properties of **Lactobacillus plantarum** CECT 7315 and CECT 7316 isolated from faeces of **healthy children**.

http://bibliome.jouy.inra.fr/demo/food/alvisir/webapi/search

# A high-precision machinery

- Artificial intelligence
- Natural Language Processing
- Machine Learning





*Alvis* software components and framework

Component 1   Component 2   ...   Component N

Ontology   ML Model   Term Lexicon   Typesystem

Scholarly Content

Corpus Building Process

Processed Output

Scientific papers and databases

**Pub Med**

Biotope ontology

INRA SCIENCE & IMPACT

efsa
European Food Safety Authority

Organism taxonomy

National Center for Biotechnology Information
NCBI

Bioinformatics

ifb

mi g le
Plateforme de BioInformatique – INRA Jouy en Josas

IBiSA

**Documents**   **Knowledge Resources**   **Client Application**

# OntoBiotope in an European Open Science perspective

OntoBiotope service becomes an application of OpenMinTeD text-mining infrastructure
**Benefits** from

- **Full-text paper collection aggregation, standardisation**
  - *OpenAire, CORE*

- Guaranteed **computational resources** in a secure environment, virtual machines and monitoring capabilities
  - *Okeanos service of GRNet in EGI federated cloud.*

- **Semantic resources** aggregation, uniform access, standard representation, update
  - *AgroPortal* (integration *in progress,* Visa TM project)

# OntoBiotope in an European Open Science perspective

**Deployment of OntoBiotope TDM on OpenMinTed infrastructure offers to the scientific communities**

A fully open access in a unified framework to the service, the processing workflows, the input data, the TDM final and intermediate results
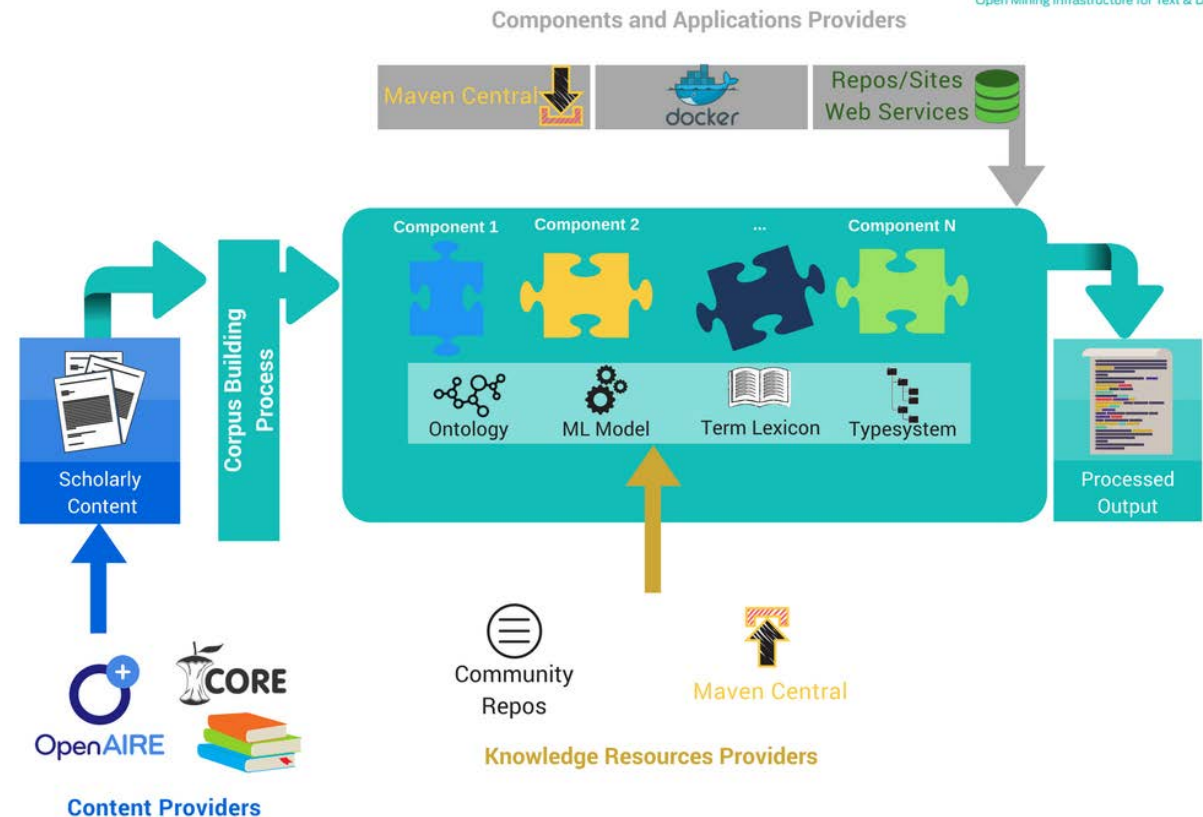
Reproducibility and adaptability.
Non TDM specialists (bioinformaticians) can reuse the workflows and replace subparts from component catalogue

The technological choices (Galaxy framework, Docker and Maven) make OMTD evolutive and interoperable.

**Innovative research**
national, European, international
**shared resources and infrastructures**

# Barriers and opportunities in the European Science Cloud

Extend document sources, remove legal and technical hindrances

1 160 000 on line articles
8 670 journals
33 plateforms

50 % accessible
to INRA researchers
OA and subscription

We do not want pdf

13 % with a text-mining clause

What food microbiology researchers need is *all* the scientific public information

- Not only the one that is Open Access
- Not only the one that is automatically findable
- Not only the one that is in a standard parsable format

# Barriers and opportunities in the European Science Cloud

**Facilitating access and use of high-level services in AgroFood.**

**In Agriculture and Life Sciences, text-mining is not the end of the story**
TDM services and results will be combined and integrated into wider data analysis applications
Inteleave text-mining and experimental data analysis
Experts may intervene in the analysis.

**Improve OMTD virtual research environment** to support collaboration and sharing data, knowledge models and workflows

**FAIR data**
Federating more data and knowledge from all fileds of food microbiology : genetics, biological resources, health, nutritional information, industrial process, distribution, retail, cooking