



FLORILEGE: an integrative database using text mining and ontologies

Estelle Chaix, Sandra Derozier, Louise Deleger, Hélène Falentin,
Jean-Baptiste Bohuon, Mouhamadou Ba, Robert R. Bossy, Delphine Sicard,
Valentin Loux, Claire Nédellec

► To cite this version:

Estelle Chaix, Sandra Derozier, Louise Deleger, Hélène Falentin, Jean-Baptiste Bohuon, et al.. FLO-RILEGE: an integrative database using text mining and ontologies. JOBIM 2018, Jul 2018, Marseille, France. , 2018, ABSTRACTS JOBIM 2018. hal-01827946

HAL Id: hal-01827946

<https://hal.science/hal-01827946>

Submitted on 27 Sep 2019

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

FLORILEGE: an integrative database using text mining and ontologies

Estelle Chaix¹, Sandra Derozier¹, Louise Deléger¹, Hélène Falentin², Jean-Baptiste Bohuon¹, Mouhamadou Ba¹, Robert Bossy¹, Delphine Sicard³, Valentin Loux¹, Claire Nédellec¹

¹ MalAGE, INRA, Université Paris-Saclay, 78350, Jouy-en-Josas, France ² STLO, INRA, Agrocampus Ouest, 35042, Rennes, France

³ SPO, INRA, Université Montpellier, 34000, Montpellier, France

Biological question : What microorganisms live in my food?

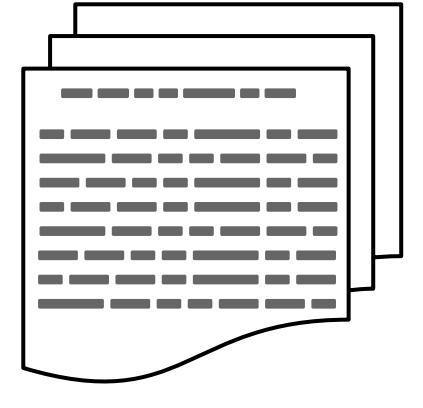
In recent years, developments in molecular technologies have led to an exponential growth of experimental data and publications spread over multiple sources. Therefore, researchers need applications, that provide an unified access to both data and related scientific articles.

The design of dedicated applications and services requires infrastructures and tools such that application developers and data managers can easily access to and process textual data, link them with other data and make the results available to scientists.

Florilege application dedicated to Food Microbiology is an example of application built on the top of the OpenMinTeD infrastructure.

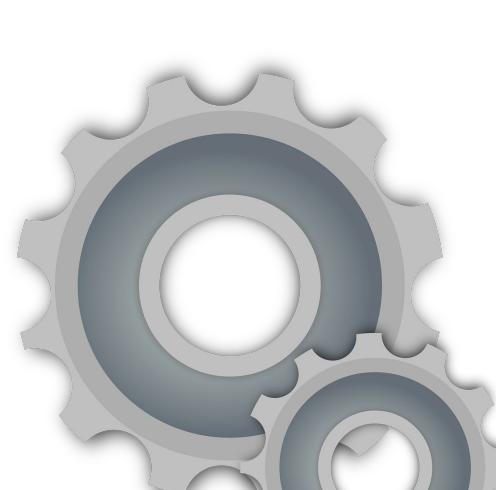
Text mining

> scientific publications?



The effect of high hydrostatic pressure on the survival of the psychrotrophic organisms *Listeria monocytogenes*, *Bacillus cereus*, and *Pseudomonas fluorescens* was investigated in ultrahigh-temperature milk.

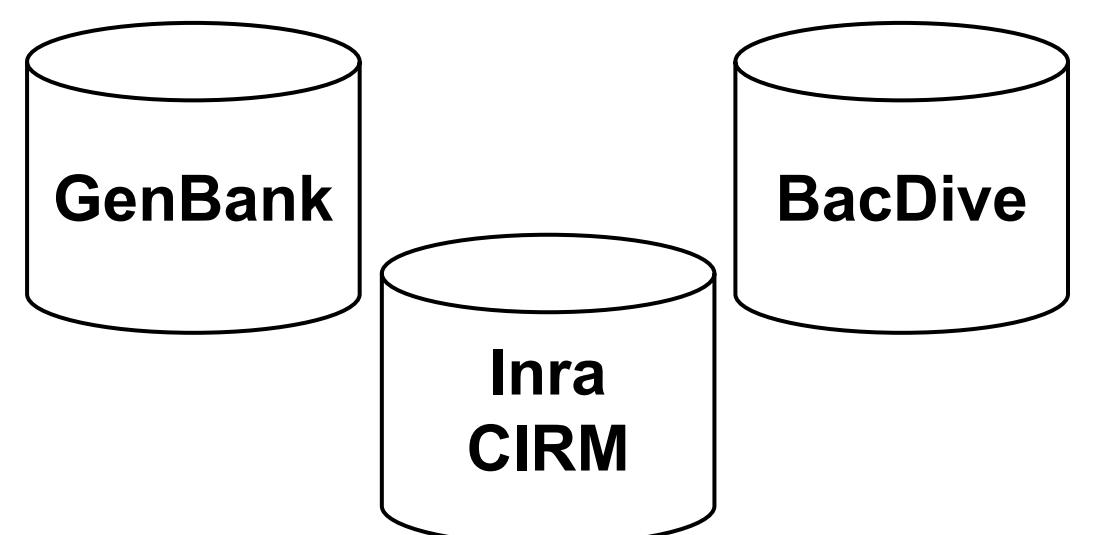
Named Entity Recognition



Detection of relevant biological entities to answer the biological question

How to process the text from:

> databases?



Knowledge resources

Taxonomies and ontologies

Formal structured representations such as ontologies provide a shared reference representation (Kelsel et al., 2010) for heterogeneous information from various sources. Ontologies also overcome the limitations of keyword-based search engines: semantic search engines extend simple string-matching with query facility on general terms that provide answers independently of how they are expressed in the searched text (Chaix et al., 2018).

The Ontobiotope ontology is in a formal machine-readable representation that enables indexing of information as well as conceptualization and reasoning.

Ontobiotope ontology is available on Agroportal : <http://agroportal.lirmm.fr/ontologies/ONTOBIOTOP>

FLORILEGE database

Aggregation of heterogeneous data

SOURCE TEXT	HABITAT	RELATION TYPE	TAXON	SOURCE
9798141	cheese	is inhabited by	<i>Escherichia</i>	OpenMinTeD
21338778	cheese	is inhabited by	<i>Bifidobacterium</i>	OpenMinTeD
9713765, 12358494, 23349056	cheese	is inhabited by	<i>Penicillium</i>	OpenMinTeD
HM462426, HM462423, AB326301	cheese	is inhabited by	<i>Lactobacillus plantarum</i>	GenBank
26082116	cheese	is inhabited by	<i>Lactobacillus delbrueckii</i>	OpenMinTeD
25017295	cheese	is inhabited by	<i>Lactobacillus helveticus</i>	OpenMinTeD
24407037	cheese	is inhabited by	<i>Penicillium rubrum</i>	OpenMinTeD
23541205, 9276789, 11375183	Habitat: cheese	Appears in the text as:	<i>Lactobacillus</i>	MinTeD
12010558			<i>Lactobacillus acetotolerans</i>	MinTeD
8573524			<i>Lactobacillus acetotolerans</i>	MinTeD
10481407, 10618286, 19901253			<i>Lactobacillus acidifarinae</i>	MinTeD
7516360			<i>Lactobacillus acidifarinae DSM 19394</i>	MinTeD
8434927, 18468028, 22154239			<i>Lactobacillus acidiplicis</i>	MinTeD
25612091			<i>Lactobacillus acidiplicis DSM 15836</i>	MinTeD
21219740			<i>Lactobacillus acidiplicis KCTC 13900</i>	MinTeD
18331739, 11168636			<i>Lactobacillus acidophilus</i>	MinTeD
18510560			<i>Lactobacillus acidophilus 30SC</i>	MinTeD
10742208			<i>Lactobacillus acidophilus ATCC 4796</i>	MinTeD
440405, 26320771, 25998659			<i>Lactobacillus acidophilus CFH</i>	MinTeD
21740724			<i>Lactobacillus acidophilus CIP 76.13</i>	MinTeD
			<i>Lactobacillus acidophilus CIRM-BIA 442</i>	MinTeD
			<i>Lactobacillus acidophilus CRBIP 24179</i>	MinTeD
			<i>Lactobacillus acidophilus DSM 20079</i>	MinTeD
			<i>Lactobacillus acidophilus DSM 20242</i>	MinTeD
			<i>Lactobacillus acidophilus DSM 9126</i>	MinTeD
			<i>Lactobacillus acidophilus JV3179</i>	MinTeD
			<i>Lactobacillus acidophilus La-14</i>	MinTeD
			<i>Lactobacillus acidophilus NCIM</i>	MinTeD

NCBI taxonomy

Detection and normalization of:

Class of entities represents:

Example: *Lactobacillaceae*

Lactobacillus

Lactobacillus acetotolerans

Lactobacillus acetotolerans

Lactobacillus acidifarinae

Lactobacillus acidifarinae DSM 19394

Lactobacillus acidiplicis

Lactobacillus acidiplicis DSM 15836

Lactobacillus acidiplicis KCTC 13900

Lactobacillus acidophilus

Lactobacillus acidophilus 30SC

Lactobacillus acidophilus ATCC 4796

Lactobacillus acidophilus CFH

Lactobacillus acidophilus CIP 76.13

Lactobacillus acidophilus CIRM-BIA 442

Lactobacillus acidophilus CRBIP 24179

Lactobacillus acidophilus DSM 20079

Lactobacillus acidophilus DSM 20242

Lactobacillus acidophilus JV3179

Lactobacillus acidophilus La-14

Lactobacillus acidophilus NCIM

NCBI taxon 1579
Lactobacillus acidophilus
synonym: *Thermobacterium intestinalis*, *Bacillus acidophilus*

Ontobiotope Habitat

(Papazian et al., 2012; Bossy et al., 2015)

Habitat

Habitats with similar physico-chemical characteristics

milk and milk product
butter
buttermilk
cheese
fermented cheese
brined cheese
fermented fresh cheese
ripened cheese
stretched curd cheese
whey cheese
fresh cheese
chhana
cottage cheese
fermented cottage cheese
fermented fresh cheese
quark
queso blanco
queso fresco
processed cheese

Ontobiotope OBT: 0000194
fermented cheese
is_a: cheese
is_a: fermented dairy product

Phenotype

Phenotypes describing the impact of a factor on microbial behaviour

- microbial phenotype
- phenotype wrt adhesion
- phenotype wrt community behaviour
- phenotype wrt growth
- phenotype wrt metabolic activity
- phenotype wrt microbial-host interaction
- phenotype wrt morphology
- phenotype wrt motility
- motile
- non motile
- taxon phenotype
- phenotype wrt ploidy
- phenotype wrt stress
- phenotype wrt chemical composition
- acid resistant
- acid sensitive
- acid tolerant
- alkali resistant

Ontobiotope OBT: 0000372
acid resistant
synonym: acidoresistant,acidresistant
is_a: phenotype wrt chemical composition
is_a: stress resistant

Entity categorization

Categorization allows to abstract and formalize the extracted entities from the form of the raw text to a generic class

Lactococcus fujisensis Lives in outer leaves of Chinese cabbages
NCBI taxon 610251 Lactococcus fujisensis Ontobiotope OBT:0000094 Chinese cabbage
L. fujisensis Lives in fermented soybean food samples
NCBI taxon 610251 Lactococcus fujisensis Ontobiotope OBT:0000433 soybean and related product
+ Ontobiotope OBT:0001668 fermented soybean
Listeria monocytogenes Exhibits psychrotrophic
NCBI taxon 1639 Listeria monocytogenes Ontobiotope OBT: 0000328 psychrotrophic

Perspectives

- ★ Web application available at <http://migale.jouy.inra.fr/Florilege/>
- ★ A database with predicted relationships by text mining
 - Taxon ↔ Habitat (820,000 relations)
 - Taxon ↔ Phenotype (86,000 relations)
- ★ A direct link to external public data (DSMZ, GenBank, CIRM)
- ★ Hierarchical and synonym search
- ★ Query filtering on data source, QPS status...
- ★ Data export in tabulated format

Acknowledgements

This work was supported by the OpenMinTeD project (EC/H2020-EINFRA 654021). We would like to thank the biologists of the Florilège working group of the metaprogramme MEM - Meta-omics and microbial ecosystems- of the French National Institute for Agricultural Research (Inra) and the Food Microbiome project, for their participation in the enrichment of the OntoBiotope Habitat ontology.

