



FLORILEGE: an integrative database using text mining and ontologies

Estelle Chaix, Sandra Derozier, Louise Deleger, H  l  ne Falentin,
Jean-Baptiste Bohuon, Mouhamadou Ba, Robert R. Bossy, Delphine Sicard,
Valentin Loux, Claire N  dellec

► To cite this version:

Estelle Chaix, Sandra Derozier, Louise Deleger, H  l  ne 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 destin  e au d  p  t et    la diffusion de documents scientifiques de niveau recherche, publi  s ou non,   manant des   tablissements d'enseignement et de recherche fran  ais ou   trangers, des laboratoires publics ou priv  s.

Biological question : What microorganisms live in my food?

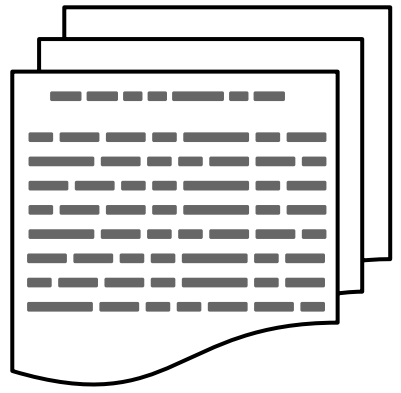
In recent years, developments in molecular technologies have led to an exponential growth of experimental data and publications spread over multiple sources. Therefore, researchers need applications, that provide an unified access to both data and related scientific articles.

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.

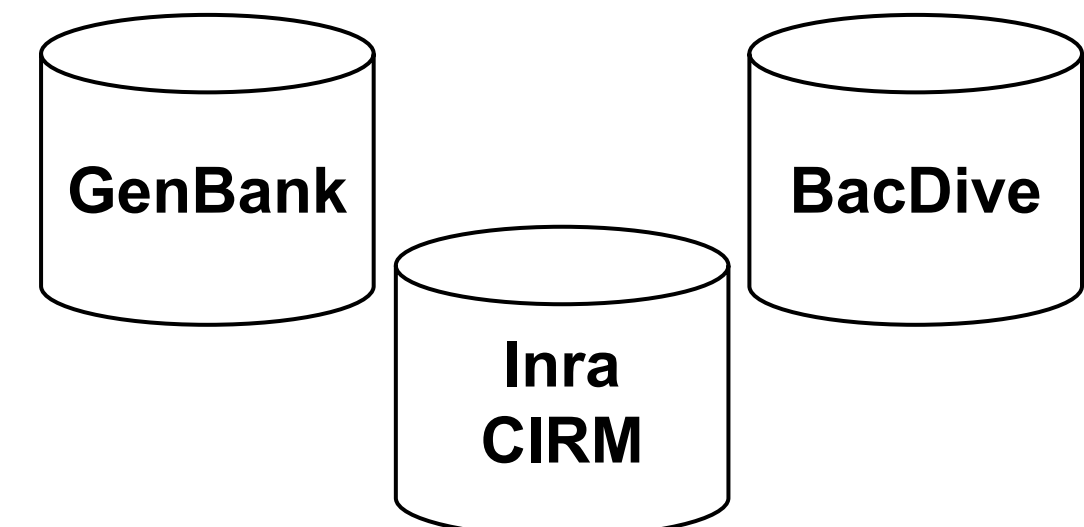
Text mining

➤ 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.

➤ databases?



Named Entity Recognition



Detection of relevant biological entities to answer the biological question

What information?

★ **Entities** (terms with particular interest for microbiologists)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

★ **Relationships** (links between entities)

Critical information stored in

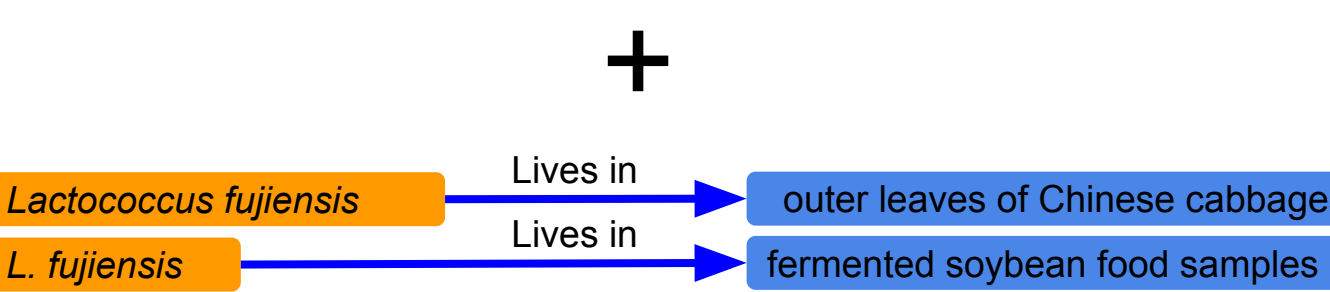
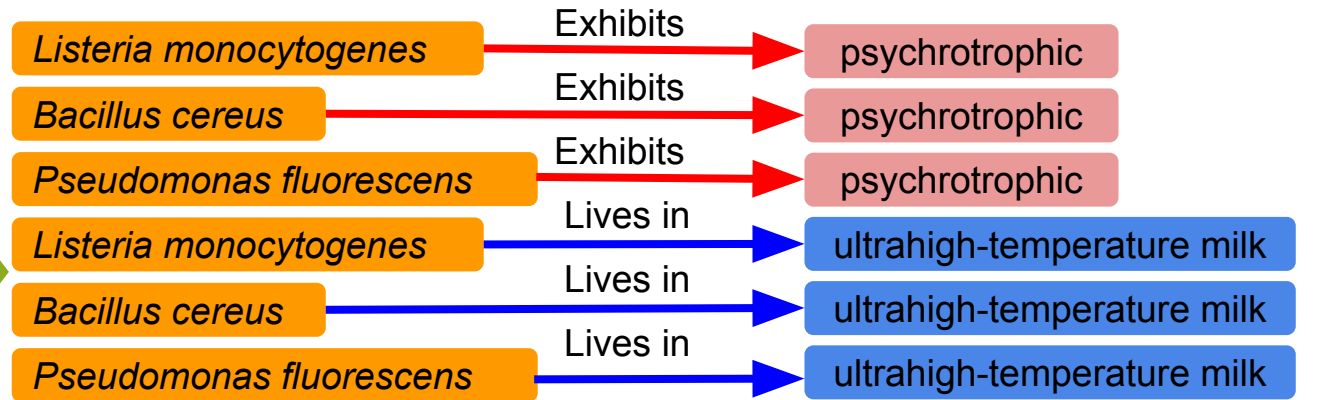
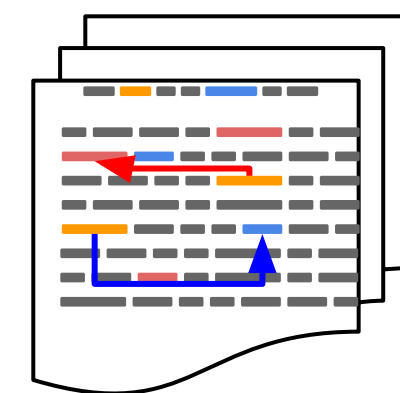
★ public databases

- GenBank (NCBI)
- GOLD (JGI)
- BacDive (DSMZ)
- CIRM database (INRA)

★ scientific publications

- Abstract (e.g. PubMed)
- Open access full texts (e.g. PMC)
- Non-open access publications

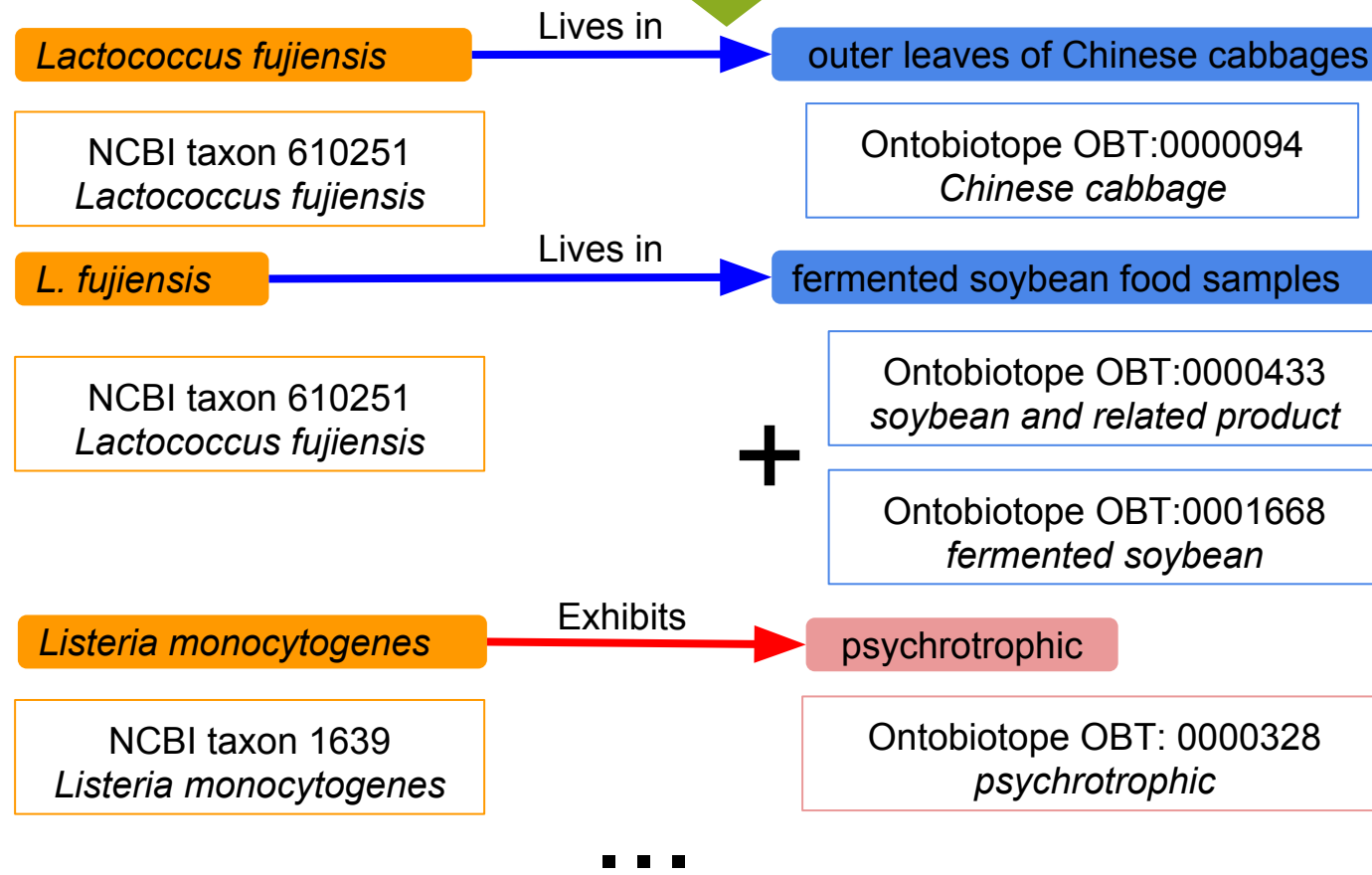
This workflow is available on the OpenMinTeD platform : Alvis Habitat-Phenotype Relation Extractor for Microbes
<https://frama.link/FlorilegeWorkflow>



A key step for heterogeneous data integration

Entity categorization

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



Knowledge resources

Taxonomies and ontologies

Formal structured representations such as ontologies provide a shared reference representation (Kelso *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/ONTOBIOTOPE>

FLORILEGE database



Aggregation of heterogeneous data

| Source Text | Habitat | Relation Type | Taxon | Source |
|------------------------------|-----------------|-----------------|--|------------|
| 9795141 | cheese | is inhabited by | Escherichia | OpenMinTeD |
| 21338778 | cheese | is inhabited by | Bifidobacterium animalis | OpenMinTeD |
| 9713785, 12358494, 23349056 | cheese | is inhabited by | Penicillium raistrickii | OpenMinTeD |
| HM462426, HM462423, AB326301 | cheese | is inhabited by | Lactobacillus plantarum | GenBank |
| 26082116 | cheese | is inhabited by | Lactobacillus delbrueckii | OpenMinTeD |
| 25017295 | cheese | is inhabited by | Lactobacillus helveticus | OpenMinTeD |
| 24407037 | cheese | is inhabited by | Penicillium rubens | OpenMinTeD |
| 23541205, 9276789, 11375183 | Habitat: cheese | is inhabited by | Taxon: Lactobacillus | MinTeD |
| 12010558 | cheese | is inhabited by | Lactobacillus delbrueckii subsp. lactis, Lactobacillus paraplantarum, Lactobacillus acidophilus NCIM, Lactobacillus delbrueckii, Lactobacillus plantarum, Lactobacillus brevis, Lactobacillus mucosae, Lactobacillus rhamnosus, Lactobacillus sp. RY2, Lactobacillus helveticus DPC4571, Lactobacillus casei group, Lactobacillus paracasei, Lactobacillus parabuchneri, Lactobacillus plantarum WVE 92, Lactobacillus, Lactobacillus casei subsp. casei, Lactobacillus helveticus, Lactobacillus paracasei subsp. paracasei BGG-22-2, Lactobacillus wasatchensis, Lactobacillus coryniformis, Lactobacillus harbinensis, Lactobacillus sakei, Lactobacillus helveticus DPC 4571, Lactobacillus acidophilus, Lactobacillus, Lactobacillus reuteri, Lactobacillus buchneri, Lactobacillus gasseri, Lactobacillus fermentum, Lactobacillus pentos, Lactobacillus acidipiscis, Lactobacillus gasseri K7, Lactobacillus dolivorans, Lactobacillus salivarius, Lactobacillus sp., Lactobacillus rhamnosus GG, Lactobacillus sakei subsp. sakei, Lactobacillus | MinTeD |
| 18331739, 11168536 | cheese | is inhabited by | Lactobacillus delbrueckii subsp. lactis, Lactobacillus paraplantarum, Lactobacillus acidophilus NCIM, Lactobacillus delbrueckii, Lactobacillus plantarum, Lactobacillus brevis, Lactobacillus mucosae, Lactobacillus rhamnosus, Lactobacillus sp. RY2, Lactobacillus helveticus DPC4571, Lactobacillus casei group, Lactobacillus paracasei, Lactobacillus parabuchneri, Lactobacillus plantarum WVE 92, Lactobacillus, Lactobacillus casei subsp. casei, Lactobacillus helveticus, Lactobacillus paracasei subsp. paracasei BGG-22-2, Lactobacillus wasatchensis, Lactobacillus coryniformis, Lactobacillus harbinensis, Lactobacillus sakei, Lactobacillus helveticus DPC 4571, Lactobacillus acidophilus, Lactobacillus, Lactobacillus reuteri, Lactobacillus buchneri, Lactobacillus gasseri, Lactobacillus fermentum, Lactobacillus pentos, Lactobacillus acidipiscis, Lactobacillus gasseri K7, Lactobacillus dolivorans, Lactobacillus salivarius, Lactobacillus sp., Lactobacillus rhamnosus GG, Lactobacillus sakei subsp. sakei, Lactobacillus | MinTeD |
| 18910260 | cheese | is inhabited by | Lactobacillus delbrueckii subsp. lactis, Lactobacillus paraplantarum, Lactobacillus acidophilus NCIM, Lactobacillus delbrueckii, Lactobacillus plantarum, Lactobacillus brevis, Lactobacillus mucosae, Lactobacillus rhamnosus, Lactobacillus sp. RY2, Lactobacillus helveticus DPC4571, Lactobacillus casei group, Lactobacillus paracasei, Lactobacillus parabuchneri, Lactobacillus plantarum WVE 92, Lactobacillus, Lactobacillus casei subsp. casei, Lactobacillus helveticus, Lactobacillus paracasei subsp. paracasei BGG-22-2, Lactobacillus wasatchensis, Lactobacillus coryniformis, Lactobacillus harbinensis, Lactobacillus sakei, Lactobacillus helveticus DPC 4571, Lactobacillus acidophilus, Lactobacillus, Lactobacillus reuteri, Lactobacillus buchneri, Lactobacillus gasseri, Lactobacillus fermentum, Lactobacillus pentos, Lactobacillus acidipiscis, Lactobacillus gasseri K7, Lactobacillus dolivorans, Lactobacillus salivarius, Lactobacillus sp., Lactobacillus rhamnosus GG, Lactobacillus sakei subsp. sakei, Lactobacillus | MinTeD |
| 10742208 | cheese | is inhabited by | Lactobacillus delbrueckii subsp. lactis, Lactobacillus paraplantarum, Lactobacillus acidophilus NCIM, Lactobacillus delbrueckii, Lactobacillus plantarum, Lactobacillus brevis, Lactobacillus mucosae, Lactobacillus rhamnosus, Lactobacillus sp. RY2, Lactobacillus helveticus DPC4571, Lactobacillus casei group, Lactobacillus paracasei, Lactobacillus parabuchneri, Lactobacillus plantarum WVE 92, Lactobacillus, Lactobacillus casei subsp. casei, Lactobacillus helveticus, Lactobacillus paracasei subsp. paracasei BGG-22-2, Lactobacillus wasatchensis, Lactobacillus coryniformis, Lactobacillus harbinensis, Lactobacillus sakei, Lactobacillus helveticus DPC 4571, Lactobacillus acidophilus, Lactobacillus, Lactobacillus reuteri, Lactobacillus buchneri, Lactobacillus gasseri, Lactobacillus fermentum, Lactobacillus pentos, Lactobacillus acidipiscis, Lactobacillus gasseri K7, Lactobacillus dolivorans, Lactobacillus salivarius, Lactobacillus sp., Lactobacillus rhamnosus GG, Lactobacillus sakei subsp. sakei, Lactobacillus | MinTeD |
| 440405, 26320771, 25998659 | cheese | is inhabited by | Lactobacillus delbrueckii subsp. lactis, Lactobacillus paraplantarum, Lactobacillus acidophilus NCIM, Lactobacillus delbrueckii, Lactobacillus plantarum, Lactobacillus brevis, Lactobacillus mucosae, Lactobacillus rhamnosus, Lactobacillus sp. RY2, Lactobacillus helveticus DPC4571, Lactobacillus casei group, Lactobacillus paracasei, Lactobacillus parabuchneri, Lactobacillus plantarum WVE 92, Lactobacillus, Lactobacillus casei subsp. casei, Lactobacillus helveticus, Lactobacillus paracasei subsp. paracasei BGG-22-2, Lactobacillus wasatchensis, Lactobacillus coryniformis, Lactobacillus harbinensis, Lactobacillus sakei, Lactobacillus helveticus DPC 4571, Lactobacillus acidophilus, Lactobacillus, Lactobacillus reuteri, Lactobacillus buchneri, Lactobacillus gasseri, Lactobacillus fermentum, Lactobacillus pentos, Lactobacillus acidipiscis, Lactobacillus gasseri K7, Lactobacillus dolivorans, Lactobacillus salivarius, Lactobacillus sp., Lactobacillus rhamnosus GG, Lactobacillus sakei subsp. sakei, Lactobacillus | MinTeD |
| 2140724 | cheese | is inhabited by | Lactobacillus delbrueckii subsp. lactis, Lactobacillus paraplantarum, Lactobacillus acidophilus NCIM, Lactobacillus delbrueckii, Lactobacillus plantarum, Lactobacillus brevis, Lactobacillus mucosae, Lactobacillus rhamnosus, Lactobacillus sp. RY2, Lactobacillus helveticus DPC4571, Lactobacillus casei group, Lactobacillus paracasei, Lactobacillus parabuchneri, Lactobacillus plantarum WVE 92, Lactobacillus, Lactobacillus casei subsp. casei, Lactobacillus helveticus, Lactobacillus paracasei subsp. paracasei BGG-22-2, Lactobacillus wasatchensis, Lactobacillus coryniformis, Lactobacillus harbinensis, Lactobacillus sakei, Lactobacillus helveticus DPC 4571, Lactobacillus acidophilus, Lactobacillus, Lactobacillus reuteri, Lactobacillus buchneri, Lactobacillus gasseri, Lactobacillus fermentum, Lactobacillus pentos, Lactobacillus acidipiscis, Lactobacillus gasseri K7, Lactobacillus dolivorans, Lactobacillus salivarius, Lactobacillus sp., Lactobacillus rhamnosus GG, Lactobacillus sakei subsp. sakei, Lactobacillus | MinTeD |

References:
Bossy, R., Golik, W., Ratkovic, Z., Valsamou, D., Bessieres, P., & Nédellec, C. (2015). Overview of the gene regulation network and the bacteria biotope tasks in BioNLP'13 shared task. *BMC bioinformatics*, 16(10), S1.
Chaix, E., Déleger, L., Bossy, R., Nédellec, C. (2018). Text-mining tools for extracting information about microbial biodiversity in food. *Food Microbiology* (In Press).
Kelso, J., Hoehndorf, R., & Prüfer, K. (2010). Ontologies in biology. In *Theory and applications of ontology: Computer applications* (pp. 347-371). Springer, Dordrecht.
Papazian, F., Bossy, R., & Nédellec, C. (2012, July). AlvisAE: a collaborative Web text annotation editor for knowledge acquisition. In *Proceedings of the Sixth Linguistic Annotation Workshop* (pp. 149-152). ACL.
Nédellec, C., Bossy, R., Chaix, E., & Déleger, L. (2018). Text-mining and ontologies: new approaches to knowledge discovery of microbial diversity. *arXiv preprint arXiv:1805.04107*.

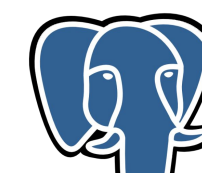
- ★ Web application available at <http://migale.jouy.inra.fr/Florilege/>
- ★ A database with predicted relationships by text mining
 - Taxon ↔ Habitat (820,000 relations)
 - Taxon ↔ Phenotype (86,000 relations)
- ★ A direct link to external public data (DSMZ, GenBank, CIRM)
- ★ Hierarchical and synonym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

Perspectives

- Update Ontobiotope on Agroportal
- Adding BioSample data & Updating GenBank data
- Taxonomy and ontology navigation on Florilege database
- Improvement of Phenotype detection by a text mining process

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.



PostgreSQL

