# Ontology-based text mining
# for microbiology research

**Claire Nédellec, MaIAGE**

**DataIA-JIST International Symposium on Data Science and AI – 10 juillet 2018**

# How can we make sense of textual data

Over 60 million articles, 2.5 per year
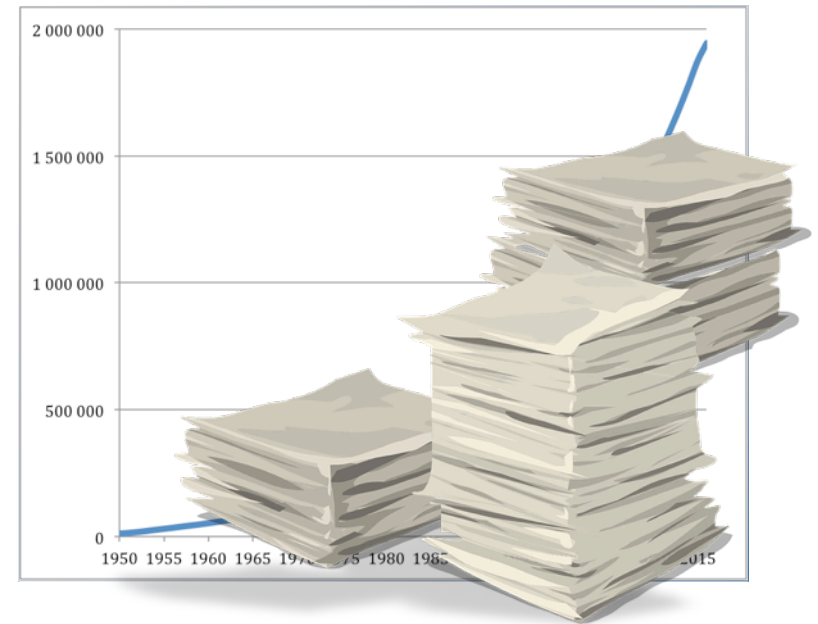160 million scientific documents indexed

50% of the papers are not read
90% of the papers are not cited
80% of the cited papers are not read
...

*The STM report, 2015*

*Orduña-Malea et al. , 2014*

*Lokman I. Meho, the rise and rise of citation analysis, 2007.*

*Simkin & Roychowdhury. Read before you cite!*
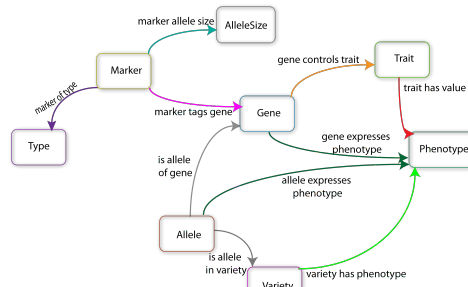




**Text-mining (TDM)**
Make sense of textual data
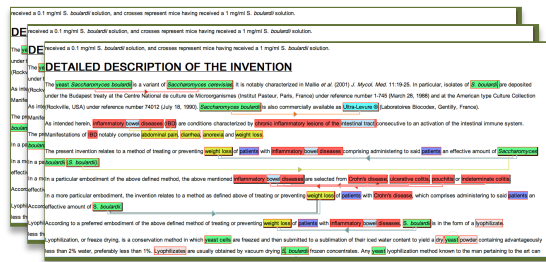Transform unstructured data into structured, machine-readable data

at the heart of the activity of non specialist researchers

# Research needs data integration, management, reasoning provided by *formal representation*
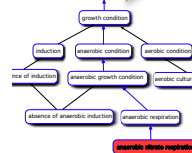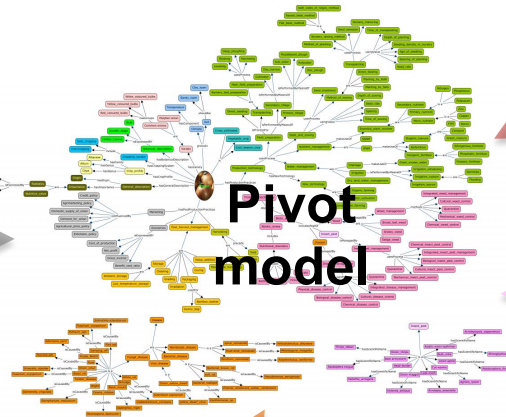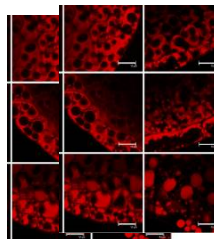


Dynamic model

Knowledge model of the text

Pivot model

Data model

Text

Image

Data

| Country | Unique Audience (000) | Time per Person (hh:mm:ss) |
|---|---|---|
| United States | 142,052 | 6:09:13 |
| Japan | 46,558 | 2:50:21 |
| Brazil | 31,345 | 4:33:10 |
| United Kingdom | 29,129 | 6:07:54 |
| Germany | 28,057 | 4:11:45 |
| France | 26,786 | 4:04:39 |
| Spain | 19,456 | 5:30:55 |
| Italy | 18,256 | 6:00:07 |
| Australia | 9,895 | 6:52:28 |
| Switzerland | 2,451 | 3:54:34 |

Source: The Nielsen Company

# Text: a source of information that requires specific treatments

Massive, diverse and under-exploited scientific information
Raising specific questions of access, analysis and interpretation

Handled at INRA by automatic *text-mining* methods
deployed on the new European infrastructure **OpenMinTeD**
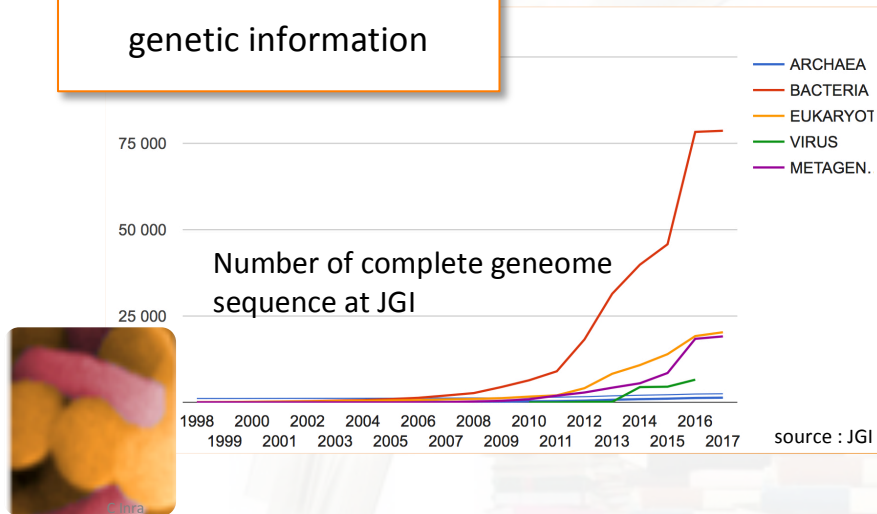
An example
- microbial biodiversity

# Microorganisms, food and scientific litterature

Billions of microorganisms everywhere, mostly unknown.

Play a critical role in food quality and transformation and its effect on health.

Microbiology research study their ecosystem and genetics for a better understanding, control and use.

Ecosystems, habitats, properties
in millions of documents

Exponential growth of
genetic information



Number of complete geneome sequence at JGI

source : JGI

Legend:
- ARCHAEA
- BACTERIA
- EUKARYOT
- VIRUS
- METAGEN.

... and publications

Number of articles about "bacteria" in PubMed bibliographic database

Pixabav

http://cm1douzant.blogspot.fr/2014/11/conte-des-droits-des-enfants-2.html

INRA
SCIENCE & IMPACT

# What microbes in my cheese?



**Inocculated microorganisms**

**"House" microbiota**

Starters
Ripening cultures

Human
Animal
Milk
Water
Air

Selection by the process

milk
pasteurisation
salting
cooking
washing

Brine
Shelf
Instruments
Cellar

Selection by the conditions

pH
Water activity
Osmolarity
Temperature

Caillé

lactic bacteria
Yeast/mould
Actinobacteria
Staphylococci

HALAB
Gammaproteobacteria
Actinobacteria
Yeast/mould

**Cheese microbiota**

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

INRA
SCIENCE & IMPACT

# Identification of microorganisms by their DNA

DNA sequence of cheese sample microorganisms

DNA        fragments

Alignment with known reference genomes

Copyright Teach the Microbiome

identified strains

Text mining

Confirm Interpret Explain

2.3 millions scientific papers

Habitat classes from ontology

Extracts of scientific articles

two L. monocytogenes cheese dairy isolates

L. monocytogenes persisting in a cold-smoked fish processing plant.

Listeria monocytogenes contamination in Chinese beef processing plants.

Listeria monocytogenes isolated from artisanal Portuguese cheses-making dairy.

the presence of L. monocytogenes in samples collected from crab processing plant
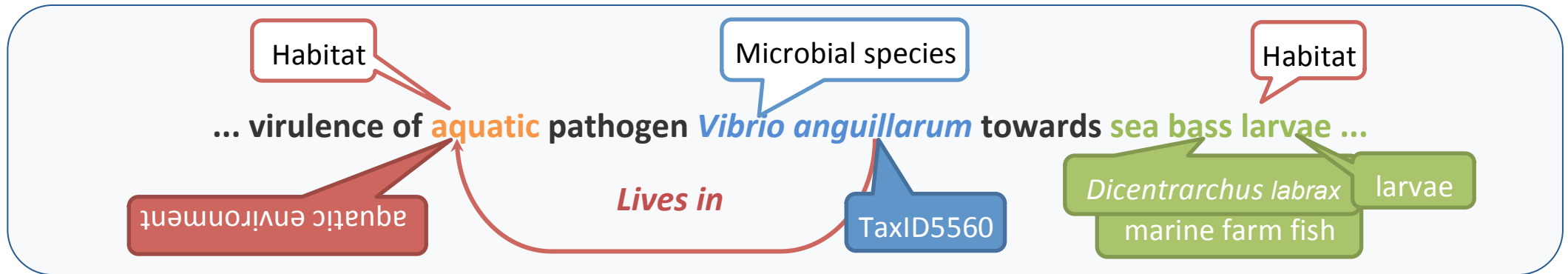
Factory

Processing factory

Food processing factory

Dairy

**TDM challenge:**
transform millions of text extracts into formal information despite the high variability

**10,000 habitats of *Listeria monocytogenes* described in PubMed**

# Information extraction



... virulence of **aquatic** pathogen *Vibrio anguillarum* towards sea bass larvae ...

Habitat — aquatic environment

Microbial species — TaxID5560 — *Lives in*

Habitat — *Dicentrarchus labrax* marine farm fish — larvae

**1. Entity recognition**    = detection (text boundaries) and broad type assignment

**2. Entity normalization**    = assignment to a category from a large set, >2,500 in OntoBiotope ontology

**3. Relationship prediction**  = links entities together over sentences, microorganism to their properties

**Artificial Intelligence methods**

Natural language processing for semantics analysis

Machine learning for generalization from examples

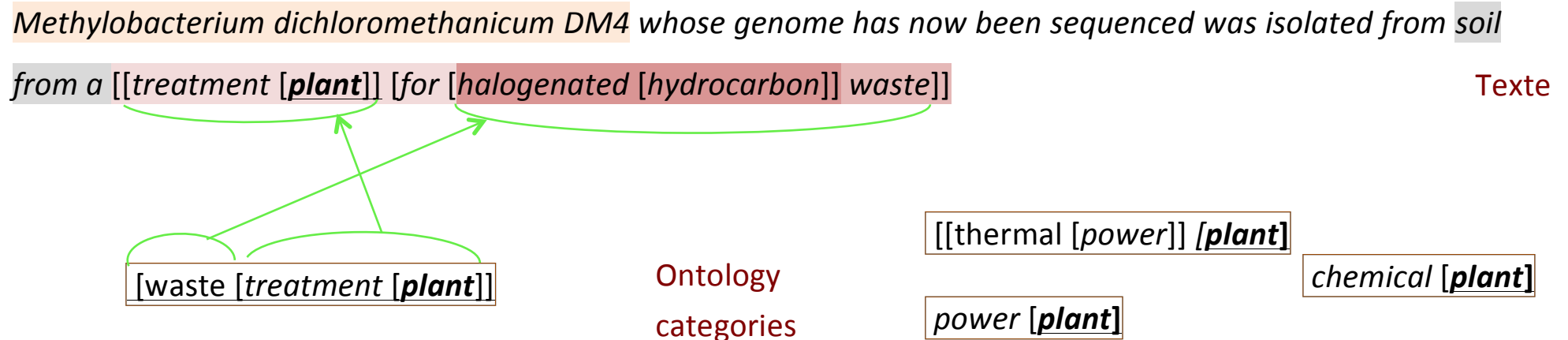Knowledge-intensive approach to deal with sparse small data

INRA SCIENCE & IMPACT

# Entity recognition and normalisation

**Entity recognition** by term extraction by BioYateA

**Entity normalization**
by *Honor*, a 2-steps method to map the text terms to the labels of the ontology categories
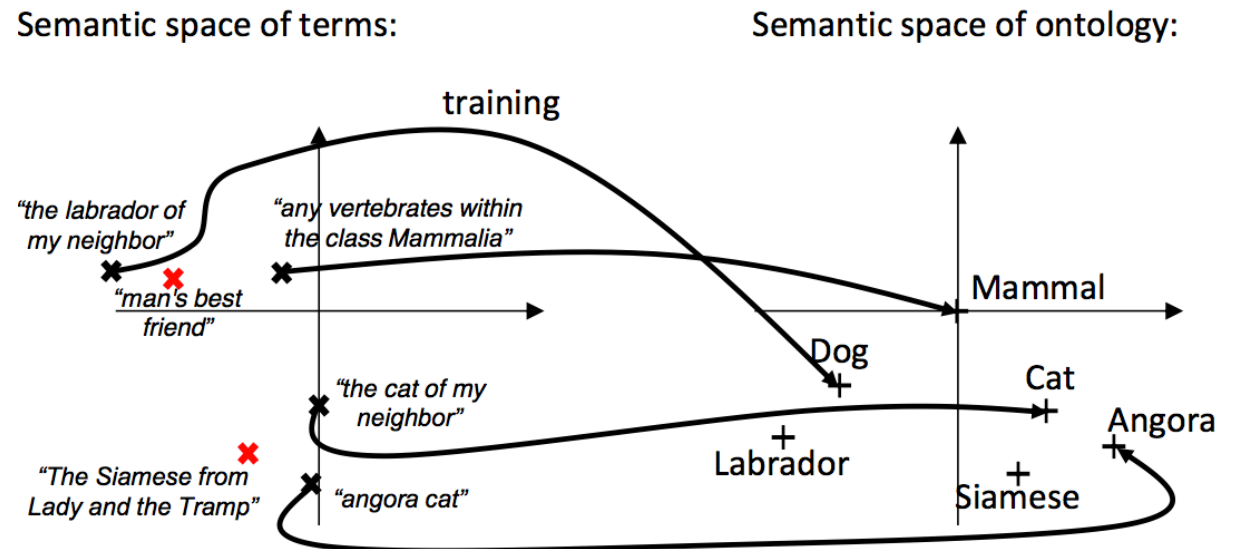
Step 1. *ToMap* computes term similarity using syntactic structures (syntactic heads and subterms)
and word similarities

*Methylobacterium dichloromethanicum DM4* whose genome has now been sequenced was isolated from soil

from a [[*treatment* [***plant***]]] [*for* [*halogenated* [*hydrocarbon*]] *waste*]]     Texte

[*waste* [*treatment* [***plant***]]]

Ontology
categories

[[thermal [*power*]] [***plant***]]

*power* [***plant***]

*chemical* [***plant***]

Ratkovic et al., BMC Bioinformatics, 2012
Golik et al., CiCLING, 2013

INRA
SCIENCE & IMPACT

# Entity recognition and normalisation

Step 2. *Contes* learns the projection of
the term vector space to the ontology
category vector space
by linear regression
from training examples
using the ontology structure

**Semantic space of terms:**

**Semantic space of ontology:**

training

"the labrador of my neighbor"

"any vertebrates within the class Mammalia"

"man's best friend"

"the cat of my neighbor"

"The Siamese from Lady and the Tramp"

"angora cat"

Mammal

Dog

Labrador

Cat

Angora

Siamese

| System | Score |
|---|---|
| HONOR with domain specific heuristics | **0.73** |
| Distant supervised HONOR | **0.72** |
| ToMap with domain specific heuristics | 0,66 |
| Turku (2017) | 0.63 |
| BOUN (2016) | 0.62 |
| ToMap | 0,61 |
| CONTES | 0.61 |

*Honor* system over performs other state of the art systems, as measured on Bacteria Biotope BioNLP-Shared Task 2016

Deléger et al., *BioNLP*, 2016

Ferré et al., *BioNLP* 2017 ; Ferré et al., *LREC* 2018

INRA
SCIENCE & IMPACT

# Automatic extraction of binary relationships, AlvisRE

| System | PRE | REC | F-M |
|---|---|---|---|
| AlvisRE | 0.51 | **0.70** | **0.59** |
| Boun revised | 0.52 | 0.53 | 0.53 |
| LIMSI revised | 0.42 | 0.60 | 0.49 |
| TEES-2.1-2 | **0.82** | 0.28 | 0.42 |
| IRISA-TexMex | 0.46 | 0.36 | 0.40 |
| Boun | 0.38 | 0.21 | 0.27 |
| LIMSI | 0.19 | 0.04 | 0.06 |

BioNLP '13: Bacteria Biotopes - Task 2

**Machine learning method based on *shortest dependency path kernel***

- Dependency path computed by CCG and abstracted by Alvis Grammar.
- Anaphora resolution
- Word distance based on word embeddings
- Global alignments computed between dependency pathes (edit distance allowing gaps + Needleman-Wunsch dynamic programming algorithm)
- Empirical Kernel Map transformation

Bossy et al., BMC Bioinformatics, 2015
Ratkovic, PhD thesis, 2014
Valsamou. PhD thesis. 2017

# Implementation and use

- All methods implemented in interoperable tools,

- combined into TDM workflows on **OpenMinTeD platform**

- Services on microbiology available at **IFP Migale Bioinformatics** plaform

  o Semantic search engine **AlvisIR**

  o **Florilege** database

- *OntoBiotope* ontology available on **AgroPortal**

openMIN7ED

# Back to microbiology: text-mining explains *Psychrobacter aquimaris* presence in cheese samples



Querying and parsing PubMed yields

- 2,3 millions documents
- 8,3 microorganisms
- 18,5 millions habitats and phenotypes assigned to more than 2500 hierarchical classes
- 7,4 millions relationships

The researcher understands :
added salt brings *P. aquimaris* to cheese

# Using Florilege database for food research and innovation

what bacteria for a new salted cheese?