

# Fouille de textes

## Applications aux données environnementales

Mathieu Roche

Cirad – TETIS – Montpellier, France



**Web :** <http://www.textmining.biz>  
**Email :** [mathieu.roche@cirad.fr](mailto:mathieu.roche@cirad.fr)



# Data Science

**Volume**

**Velocity**

**Variety**

***3V of Big Data***



**Variability, Veracity, Value,  
Visualisation, Valorization**

→ **Pluridisciplinary domain**



# Data Science

**Volume**

**Velocity**

***3V of Big Data***

**Variety**



**Variability, Veracity, Value,  
Visualisation, Valorization**

→ **Pluridisciplinary domain**



# Textual Data Science

```

0/1 x B_1404 [WARNING]: "Asynchronous reset/set/load <%item> exists in module/unit"
0/1 x B_1405 [WARNING]: "<%value> asynchronous resets in this unit detected"
0/1 x B_1406 [WARNING]: "<%value> synchronous resets in this unit detected"
0/1 x B_1407 [ERROR]: "Do not use active high asynchronous reset/set/load"

```

// Total Module Instance Coverage Summary

lines  
statements

Policy:  
<vi  
---  
0/1

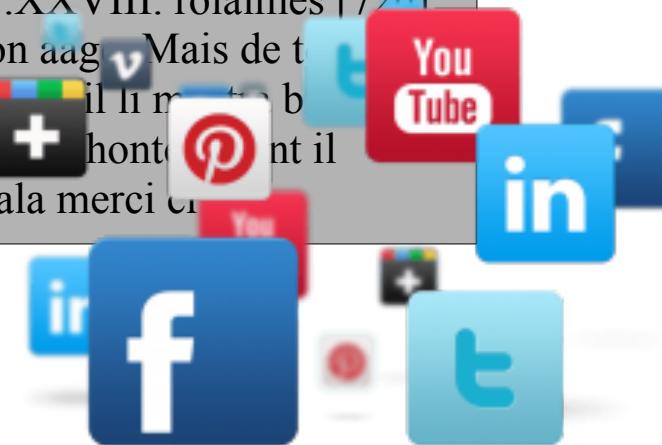


PERCENT

31.54  
31.54

>]:<message>  
-----  
is not allowed to be used as

descovri son corage a Lancelot et dist  
a guerre commença, baoit il a tot le  
e: et bien i parut, kar il fu a vint et cinc  
puis conquist il .XXVIII. roialmes [72<sup>d</sup>]  
ns fu la fin de son aag Mais de t  
st Lancelos arie, il li m'ea b  
grant honor sa + honteant il  
e roi Artu et il li ala merci or



# Pluridiciplinary projects

## Textual data and satellite images – ANIMITEX (2013-2014)

### Vakinankaratra – L'agriculture de conservation lancée

17.12.2014 | 7:18 | Non classé | 0

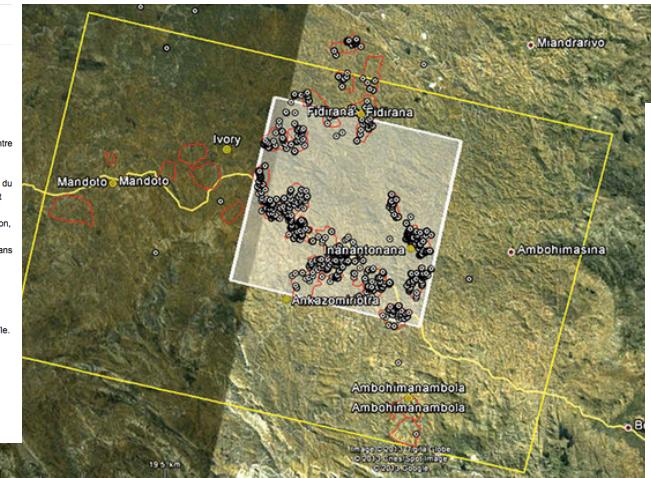


L'agro-écologie est une nécessité. Plus de 80% de la population malgache vit en milieu rural et opère en général dans l'agriculture. La croissance démographique associée au changement climatique provoque une forte destruction de l'environnement et une dégradation alarmante de la fertilité des sols. Afin d'y faire face et pour mieux lutter contre la malnutrition, le Groupement semi direct de Madagascar lance le projet Manitra dans quatre communes rurales du district de Betafika et de Mandoto, dans la région Vakinankaratra. Ce projet est réalisé en partenariat avec le ministère de l'Agriculture et du développement rural et sur financement de l'Association française du développement et du Comex.

Le groupement qui focalise son activité sur l'agro-écologie et l'agriculture de conservation, sensibilise et incite les paysans des communes ciblées à pratiquer l'agriculture sous couverture végétale et la rotation culturelle. Et afin d'assurer une sécurité alimentaire dans la commune rurale d'Ankazomirina, d'Ivainy et de Fidirana, le projet Manitra compte adhérent 1000 paysans, dont 200 femmes, sur la pratique de ce système de culture agro-écologique qui ne nécessite pas des nombreux travaux et évènements comme l'excavation et l'arrachage. « Il suffit que les paysans recourent à des méthodes de culture et d'ensemencement simples pour avoir l'argent pour l'achat d'outils », note Rakotonirina Ranaivo, directeur exécutif du projet qui l'active aussi dans le Sud-Est de l'île.

Des formations sur la régénération de la fertilité du sol et la lutte contre sa dégradation ainsi que l'introduction du système des légumineuses seront la priorité des activités du projet.

**Angola Ny Avo**



Archive ouverte  
des publications  
du Cirad

Recherche avancée | Autres Cirad | Parcourir | Déposer | Se connecter | Crise hier, opportunités aujourd'hui, défis pour demain : le cas de la filière riz à Madagascar : [Draft]

### Les facteurs de la crise 2004-2005 sur le marché du riz

La manifestation la plus visible de la situation du marché du riz en 2004 et début 2005 est une augmentation sans précédent des prix de détail. Les causes de cette crise sont une conjonction de plusieurs facteurs, internes et externes : physiques, monétaires et politiques.

L'an dernier, les prix du riz ont normalement augmenté pendant la période de soudure, fin 2003 début 2004, mais ne sont pas redescendus en période de récolte, ils ont continué à augmenter à un rythme soutenu. La variation annuelle du prix du paddy entre récolte et soudure est habituellement de l'ordre de 50% au Lac Alaotra, elle a été de 150% en 2004-2005 [Minten et Raison, 2005]. Le prix du riz national ou importé dépassait historiquement 1000 Ar le kg entre septembre 2004 et février 2005 sur les marchés de la capitale. Si on compare l'évolution des prix du riz en 2001 et en 2004, on peut se rendre compte que les trois premiers mois de l'année il coûtaient moins cher en 2004 qu'en 2001 et 2,5 fois plus cher en novembre.

Cette hausse des prix s'est généralisée dans tout le pays. Elle s'est répercute dans l'espace : marchés urbains et ruraux, auprès de tous les agents de la filière et pour toutes les variétés de riz (vary gazy, makloka, tsipala, riz pluvial...). A titre d'exemple, dans le Moyen-Ouest, la hausse des prix du riz a été aussi importante sur les marchés situés en bord de route nationale que sur les marchés plus enclavés comme Ifanadiana (45 mm de piste en saison sèche), Vasiana (1h15mn), Mahasolo (2h30mn) ou Ambalanana (4h).

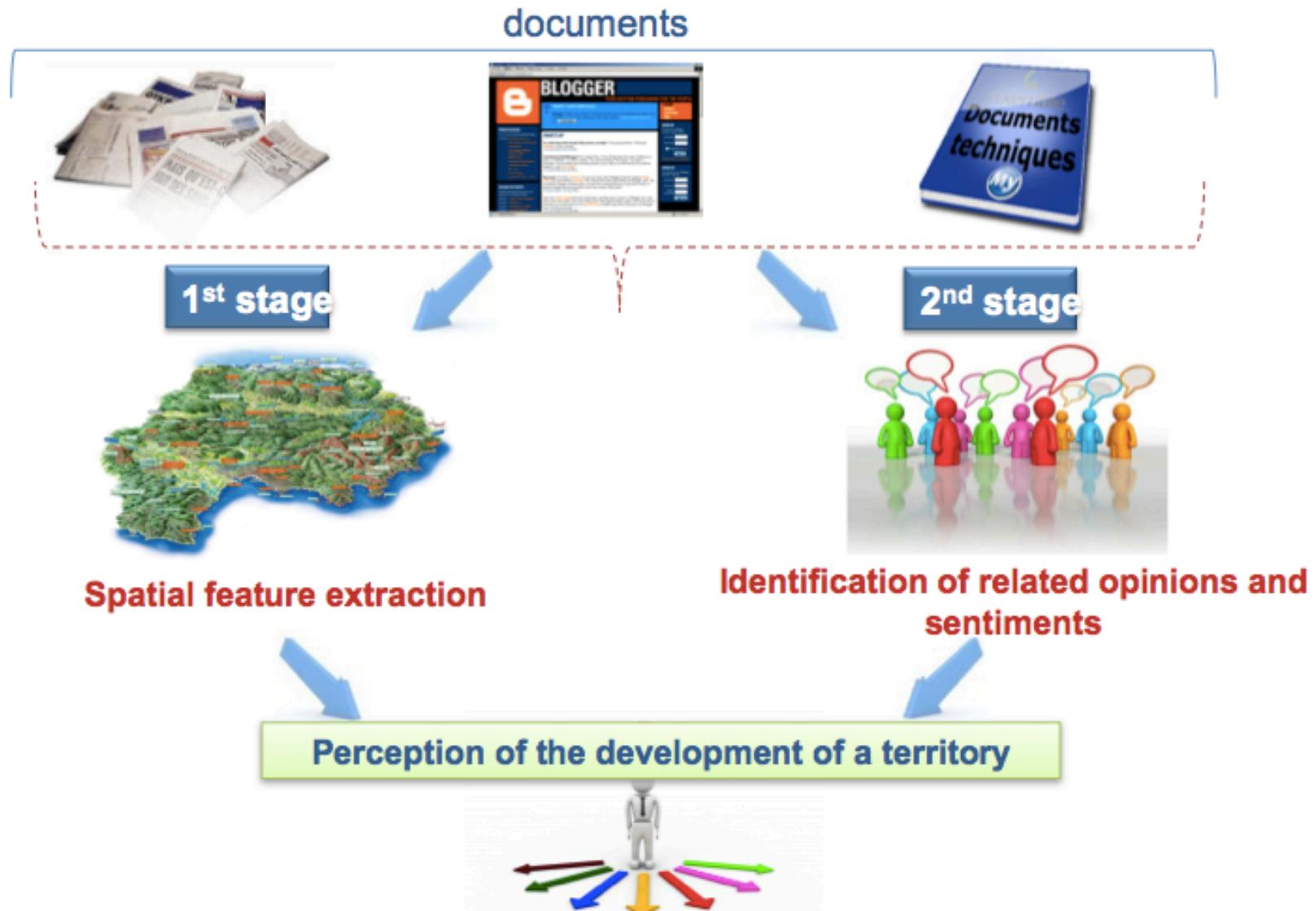
## Textual data and data bases – QuDoSSI (since 2016)



battre. Je vois un bon avenir.  
\*\*\*\* \*recit\_2095 \*sex\_Masculin \*age\_17 \*pays\_Côte-d'Ivoire  
Danané, ne se souvient pas (mai 2009), juillet 2009 C'est la situation économique de ma famille qui m'a encouragé à prendre la décision de partir. Ma famille n'était au courant de mon voyage. C'est même l'argent de mon père que j'ai volé pour entamer mon projet de voyage. C'est mon ami qui m'a parlé de cette route. Je ne peux pas vous parler de la route parce que si je le dis le chef va me punir. Le voyage s'est bien passé. Mon projet c'est d'aller jusqu'en Espagne. N.B. le mineur ne veut pas donner de détails parce qu'il a peur.  
\*\*\*\* \*recit\_2097 \*sex\_Masculin \*age\_16 \*pays\_Côte-d'Ivoire  
Lakota, ne se souvient pas, ne se souvient pas Je suis d'une famille à situation défavorable. Mon père veut qu'un enfant puisse aller en Europe et cela pouvait aider la famille. Je suis venu avec le fils d'un ami de mon père. Il a déjà fait la route. C'est lui qui me guidait. J'avais entendu parler de la route en causerie avec les amis. On est passé d'abord chez mon grand-père à Yamoussoukro. Lui aussi m'a donné de l'argent pour le voyage. C'est mon compagnon qui discutait le transport. A la demande de ma famille je lui ai donné tout mon argent. Il ne me faisait pas de compte des dépenses. C'est à cause de lui que je travaille parce qu'il a dépensé tout mon argent. Je n'ai plus d'argent pour continuer. Il est allé lui à Maghnia. Je suis resté seul ici. A Gao, on a pris le pickup jusqu'en Algérie. Nous ne sommes pas passés par la frontière officielle. Je suis rentré en Algérie sans voir la police. Mon projet migratoire c'est de rentrer au Maroc. On m'a dit qu'au Maroc je peux rentrer en Espagne même si j'ai pas d'argent. Comme je suis jeune les gens sont beaucoup gentils avec moi. Tout le monde me donne à manger. Mes employeurs me demandent mon aide pour me prendre. T1



# Senterritoire: Generic Process



# Outline

**Part 1** Data Science and Big Data

**Part 2** Textual data and heterogeneity

**Part 3** Applications in agriculture domain

**Part 4** Conclusions and future work



## Part 2

# Textual data and heterogeneity

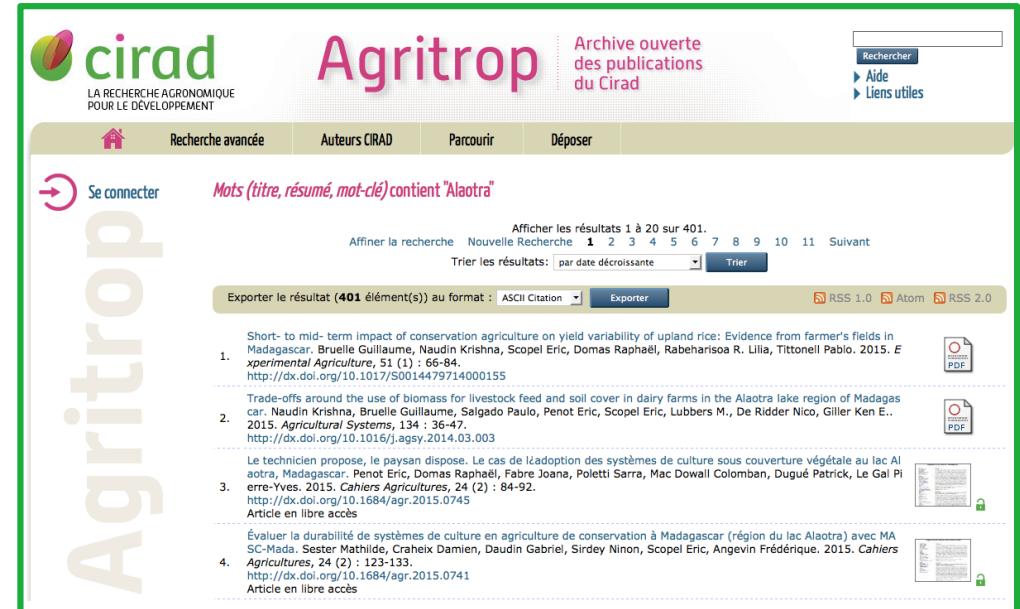


# How to match documents?

- Data and Issue**



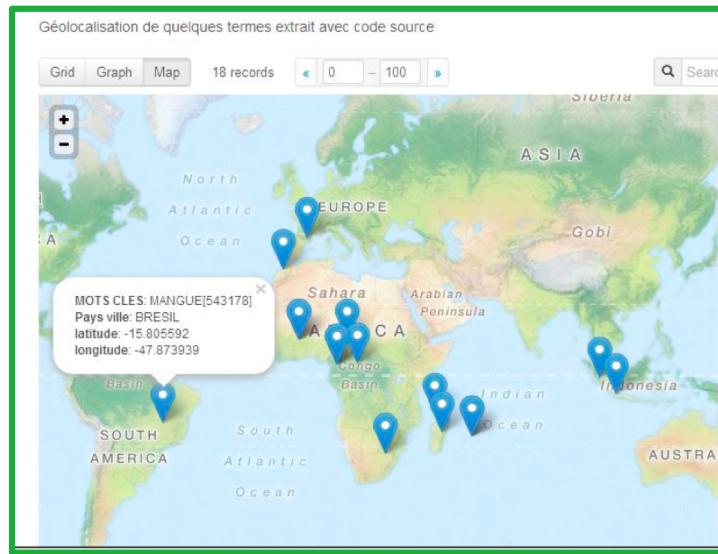
- Hard Disc (157 188 files)




- **Method: Extraction of features** [Roche et al. CA'2015]

## 3 types of features:

- thematic
- spatial
- temporal



17.12.2014 | 7:18 | Non classé | 0



L'agro-écologie est une nécessité. Plus de 80% de la population malgache vit en milieu rural et opère en général dans l'agriculture. La croissance démographique associée au changement climatique provoque une forte destruction de l'environnement et une dégradation alarmante de la fertilité des sols. Afin d'y faire face et pour mieux lutter contre la malnutrition, le Groupement semis direct de Madagascar lance le projet Manitatra dans quatre communes rurales du district de Betafo et de Mandoto, dans la région Vakinankaratra. Ce projet est réalisé en partenariat avec le ministère de l'Agriculture et du développement rural et sur financement de l'Association française du développement et du Comesa.

Le groupement qui focalise son activité sur l'agro-écologie et l'agriculture de conservation, sensibilise et incite les paysans des communes cibles à pratiquer l'agriculture sous couverture végétale et la rotation culturelle. Et afin d'assurer une sécurité alimentaire dans la commune rurale d'Ankazomirioratra, d'Inantanana, de Vinany et de Fidirana, le projet Manitatra compte adhérer 1000 paysans, dont 200 femmes, sur la pratique de ce système de culture agro-écologique qui ne nécessite pas des nombreux travaux et éreintant comme l'exige le labourage. « Il suffit que les paysans recourent le sol de végétaux et cultivent sans dépenser du temps et de l'argent pour l'achat d'outils », note Rakotondramanana, directeur exécutif du projet qui s'active aussi dans le Sud-Est de l'île. Des formations sur la régénération de la fertilité du sol et la lutte contre sa dégradation ainsi que l'introduction du système des légumineuses seront la priorité des activités du projet.



- **(a) Extraction of features: thematic terms** [Lossio Ventura et al. ISWC'2014]

## BioTex

- Système de culture
- Production
- Développement durable
- Eau ...

- Système de culture
- Développement durable
- Ressources naturelles
- Mise en œuvre ...

**Patterns Information**

Number of linguistic patterns used to filter candidate terms: 200

Patterns extracted from UMLS for English and Spanish, and from MeSH for French

Ex: Noun Noun; Noun Prep:det Noun; ... [more examples](#)

**Type of terms to extract**

All Terms       Multi Terms

single-word + multi-word term      multi-word term

**Measures selection and data**

Select ranking measure: L-value [read more](#)

Type of documents:  Single Document     Set of Documents

File source: [Parcourir...](#) Aucun fichier sélectionné.  
Only ".txt" accepted as file extension

Language of your text: English

**Extract Terms**

**Institutions**



Laboratoire  
Informatique  
Robotique  
Microélectronique  
Montpellier





**Sponsors**





- (a) Extraction of features: **spatial features (SF)**

## Model

- **Global Model:** SF is composed of at least one Named Entity (NE) and one variable number of spatial indicators specifying its location. SF can then be identified in two ways:
- **Absolute spatial feature (A\_SF)** one NE with a geo-localization, such as <(spatialIndicator)\*, NE of Location> (ex: *the city of Montpellier*).
- **Relative spatial feature (R\_SF)** one spatial with at least one SF (ex: *in the south of the city of Montpellier*).  
An R\_SF is defined as <(spatial relation)<sup>1..\*</sup>, A\_SF> or <(spatial relation)<sup>1..\*</sup>, R\_SF>  
**Five spatial relation types are considered:** orientation, distance, adjacency, inclusion, and geometric which defines union or intersection linking two SFs.



- (a) Extraction of features: spatial features (SF)

**Methods** [Kergosien *et al.*, IJGIS'2014]

- **Symbolic approach:** Using rules (*Text2Geo*) for extracting A\_SF and R\_SF

Basic patterns		Text2Geo paterns			
	A_SF	R_SF	R_SF	R_SF	OE
Precision	20%	48%	Precision	53%	84%
Recall	63%	27%	Recall	94%	66%
F-mesure	30%	34%	F-mesure	67%	74%
					50%

- **Statistic approach:** Using context and IR methods for spatial feature disambiguation [Taharat *et al.*, WIMS'2013]



- Disambiguation between **location** and **organisation**

SVM		Naive Bayes	
		SF	OE
SF		103	35
OE		44	90
<i>Accuracy</i>	<b>70.96%</b>	<i>Accuracy</i>	<b>69.12%</b>

Features with ConceptOrg		Features with ConceptSpa		Both types of features	
SF	OE	SF	OE	SF	OE
SF	108	30	SF	112	26
OE	47	87	OE	19	115
<i>Accuracy</i>	<b>71.69%</b>		<i>Accuracy</i>	<b>83.45%</b>	
<i>Accuracy</i>	<b>83.82%</b>				



- (a) Extraction of features: spatial features (SF)  
[Farvardin et al. Demo ISWC'2015]

**SENTERITTOIRE VIEW**

**DISPLAYED INFORMATION**

- Spatial Features
- Organizations
- Opinions
- Others
- Topics

**CORPUS AND DOCUMENTS**

- Corpus with different documents
  - 003\_7\_docs\_29.xml
  - CON:1\_1
  - CON:1\_2
  - CON:1\_3
  - CON:1\_4
  - CON:1\_5
  - CON:1\_6
  - CON:1\_7
  - 006\_all\_senterritoire\_100docs\_317.xml
  - 008\_all\_senterritoire\_300docs\_317.xml
  - 005\_7\_docs\_17.xml
  - 007\_all\_senterritoire\_150docs\_317.xml
  - 006\_ExtraitsTextesManuel\_Raw\_17.xml
  - 007\_multiDocs\_17.xml
  - 008\_7\_docs\_18.xml

**UPLOAD NEW CORPUS**

Upload Choose Files No file chosen

Pipeline:

French Spatial Features and Opinions

Description: This pipeline, extract all Spatial features, Organizations and Opinions from French language.

Start Process...

5

**DOCUMENTS**

A Sète, la dernière bataille perdue sur le projet de grande agglomération a refroidi, quelque peu, les ardeurs de François Commeinhes. Il vient ainsi d'opposer une fin de non recevoir au calendrier concocté par le président de Montpellier Agglomération Georges Frêche, lui demandant de ne pas délibérer le 30 avril sur l'élargissement du périmètre de Montpellier agglo à Sète et Mèze. "Je ne me laisserai pas imposer un timing par la communauté d'agglomération de Montpellier. Il ne s'agit pas de dire : "Le 18 avril (Ndlr, date du prochain conseil de l'agglomération de Thau), on y va." L'important est d'avoir un projet. Et si autour de ce projet, les discussions permettent de faire évoluer le territoire vers la CCNBT, vers Montpellier ou vers Hérault Méditerranée, pourquoi pas ? A terme, je suis persuadé que l'élargissement se fera. Mais pas dans la précipitation." C'est le message tenu vendredi dernier par le maire de Sète à Georges Frêche, au premier vice-président Jean-Pierre Moure et au directeur général des services François Delacroix, avec qui il déjeunait à Montpellier. C'est que jusque-là, le président de Montpellier agglo avait déroulé son propre timing : beaucoup plus rapide. Il comptait s'enfouir dans la brèche créée par l'échec aux municipales des principaux élus anti fusion François Liberti, Didier Sauvage, et Williams Ménez et aussi profiter de la position plus favorable de l'Etat. Montpellier agglo mettait tout en marche pour une fusion au 1er janvier 2009. L'administration de la structure intercommunale avait prévu de relancer la procédure dès le 30 avril avec le vote sur l'élargissement du périmètre de Montpellier agglo au bassin de Thau et la communauté de communes du nord du bassin de Thau. Pour calmer les inquiétudes des petites communes, le projet défendu par Montpellier agglo avait été revu à la baisse. On ne parlait plus de communauté urbaine. Le nouveau projet prévoyait la création d'une nouvelle agglomération où les communes préservaient leur droit du sol. Du coup, la nouvelle agglo se privait d'une partie de l'augmentation espérée de la dotation générale versée par l'Etat. Après le 30 avril, le préfet aurait eu deux mois pour valider le périmètre ou proposer sa propre partie. Les délégués de Sète et Mèze auraient eu de leur côté, deux autres mois pour se prononcer, à leur tour, sur le projet. Sauf que, désormais, autour du bassin de Thau, bien des élus hésitent à céder au rouleau compresseur montpelliérain. Ils sont nombreux à vouloir, d'abord, clarifier les situations financières des intercommunalités respectives avant de voir plus loin.

Des mots et des maux. Arceaux, Cévennes, Clemenceau, Pompignane, Chamberte, Beaux-Arts, Hauts-de-Massane, Saint-Martin, etc. Ils étaient tous là, hier soir - à une ou deux exceptions près, peut-être -, à avoir répondu présent à l'invitation de Patrick Vignal. Réunis à la salle Pétrarque, les représentants de ces comités ou associations de quartier ont ainsi pu, au cours de la première partie de cette réunion, mettre des mots sur leurs maux, leurs problèmes, leurs doutes ou leurs interrogations. Un florilège de paroles libérées les unes après les autres qui donne le vertige tant les situations peuvent être extrêmes, violentes, saugrenues ou opaques. Extraits. Le bien-vivre. "C'est une catastrophe, nous n'en pouvons plus de la délinquance. Dans notre quartier, il n'y a ni citoyenneté ni convivialité", s'exclamait ainsi un président de comité des Cévennes. Vite relayé par de mêmes sons de cloches du côté du Courreau et du Plan-Cabanes : "On dénonce depuis des années les histoires de drogue, les marchands de sommeil. " Et de La Pergola : "Quand quelqu'un ose ouvrir la bouche, il se retrouve le lendemain avec les pneus crevés ou la boîte à lettres cassée. " Idem pour La Pompignane et Les Aubes : "Est-ce qu'on va devoir fermer une maison pour tous à cause de toutes ces violences ? " Le manque de transparence. Du côté de Pavie, à Clemenceau, "on n'a pas assez d'info sur les comités de quartier. Pas de locaux où les trouver, de coordonnées... " À La Croix-d'Argent, aussi, on demande un local quand, dans d'autres secteurs, on déplore l'absence de dialogue et, surtout, de réponses des élus. "Je ne sais toujours pas à quoi je serai mis au courant", assène ainsi le représentant des Hauts-de-Massane. Sans parler des espoirs déçus, des conseils citoyens de secteur muets selon les quartiers. Après ces prises de paroles cathartiques, Patrick Vignal a dévoilé ses projets : conseils consultatifs de quartier, redéroulance du territoire,

3

4

**MARKE INFORMATIONS**

**Spatial Features**

search...

- du côté du Courreau et du Plan-Cabanes
- de Montpellier, Sète et Nîmes
- de Sète à Georges Frêche
- à Frêche
- au Schéma
- à Port Le Nouvel

**Organizations**

search...

- pour La Pompignane
- la Région aurait
- de Fabrèges
- Le PPRIF
- la BNP
- Montpellier agglo
- la Région de payer la CCI

**Opinions**

search...

- conflict
- déficit
- dissiper

©2015 - Senterritoire Web Application - Admin

M. Roche – Cours ECA, novembre 2016

15

- (c) Similarity

$$\text{Global\_Sim}(\text{vect1}, \text{vect2}) = \\ \alpha.\text{cosT}(\text{vect1}, \text{vect2}) + (1-\alpha).\text{cosS}(\text{vect1}, \text{vect2})$$

with  $\alpha \in [0, 1]$

$\text{cosT}$ : cosine based on **thematic features** (BioTex)

$\text{cosS}$ : cosine based on **spatial features**

**Perspective:** adding temporal information





## Part 3

# Applications in agricultural domain

## Animal disease surveillance

In collaboration with **CMAEE** lab  
(Control of exotic and emerging animal diseases)



# Why the need of Epidemic Intelligence?

More than **60% of the initial outbreak reports** come from unofficial informal and **heterogeneous sources**, including sources other than the electronic media, which **require verification** [Arsevska et al. ISVEE'2015]

**INTERNATIONAL BUSINESS TIMES**  
MONDAY, JUNE 01, 2015 AS OF 2:24 PM CDT

Home   Politics ▾   Economy ▾   Markets / Finance ▾   Companies ▾   Technology ▾

TECHNOLOGY   SCIENCE

Unknown Disease Kills Kazakhstan's Rare Saiga Antelopes, Scientists Baffled

By Kukil Bora   [@KukilBora](#)   on May 30 2015 7:21 AM EDT

**News**



African Swine Fever in Three Lithuanian Wild Boar  
18 May 2015

LITHUANIA - Three wild boar found at two locations were confirmed with African swine fever last week.



## Mysterious disease kills Nigerian patients within a day

The unknown disease has so far killed 17 people in a southeastern Nigerian town and officials have ruled out Ebola.

18 Apr 2015 21:32 GMT | [Health](#), [Nigeria](#), [Africa](#)



# How to detect disease outbreak on the Web?

- **Four animal disease models:** African swine fever (ASF), Foot-and-mouth disease (FMD), Bluetongue (BTV), and Schmallenberg virus (SBV)
- **First studying model:** ASF



[Industries](#) | Wed Jul 23, 2014 9:54am EDT

Related: NON-CYCICAL CONSUMER GOODS

## Poland investigates suspected case of African swine fever in farm pigs

WARSAW, JULY 23

Polish local authorities said on Wednesday that preliminary tests have pointed to a case of African swine fever (ASF) among farm pigs in eastern Poland near the city of Bialystok.

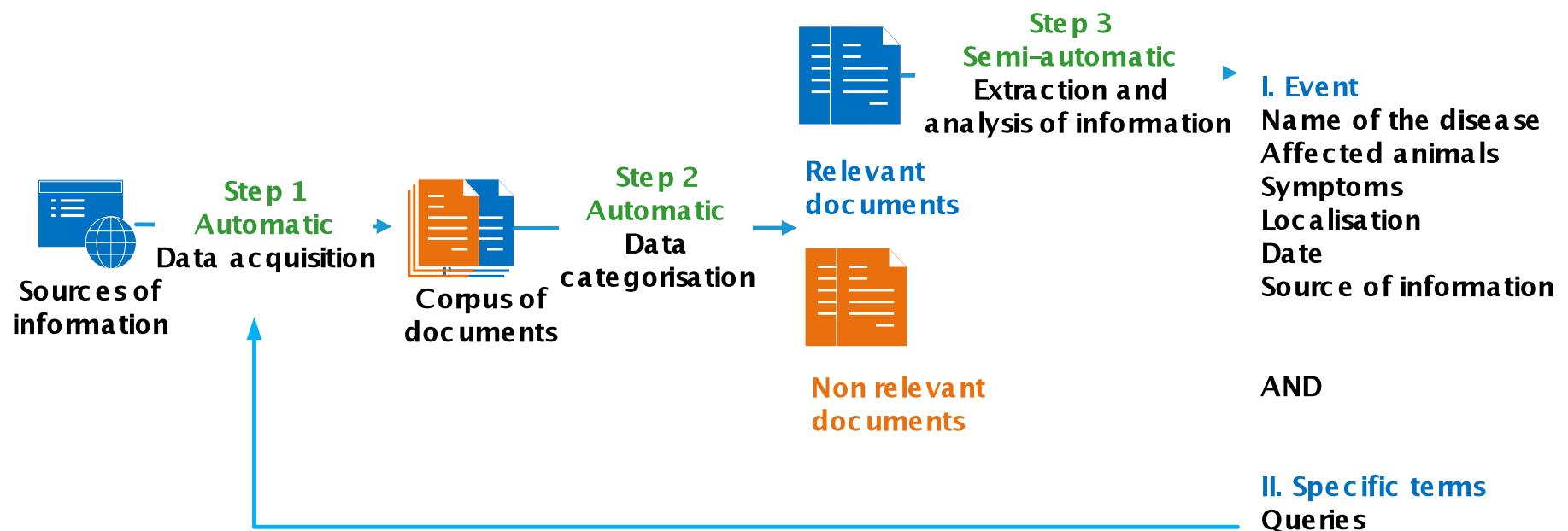
The head of the Grodziec county, Wieslaw Kulesza, told Reuters that preliminary results of tests showed that ASF was the cause of death of two-three farm pigs in the county.

"We are marking the area," Kulesza said, adding that further steps, such as laying special mats, were being taken.

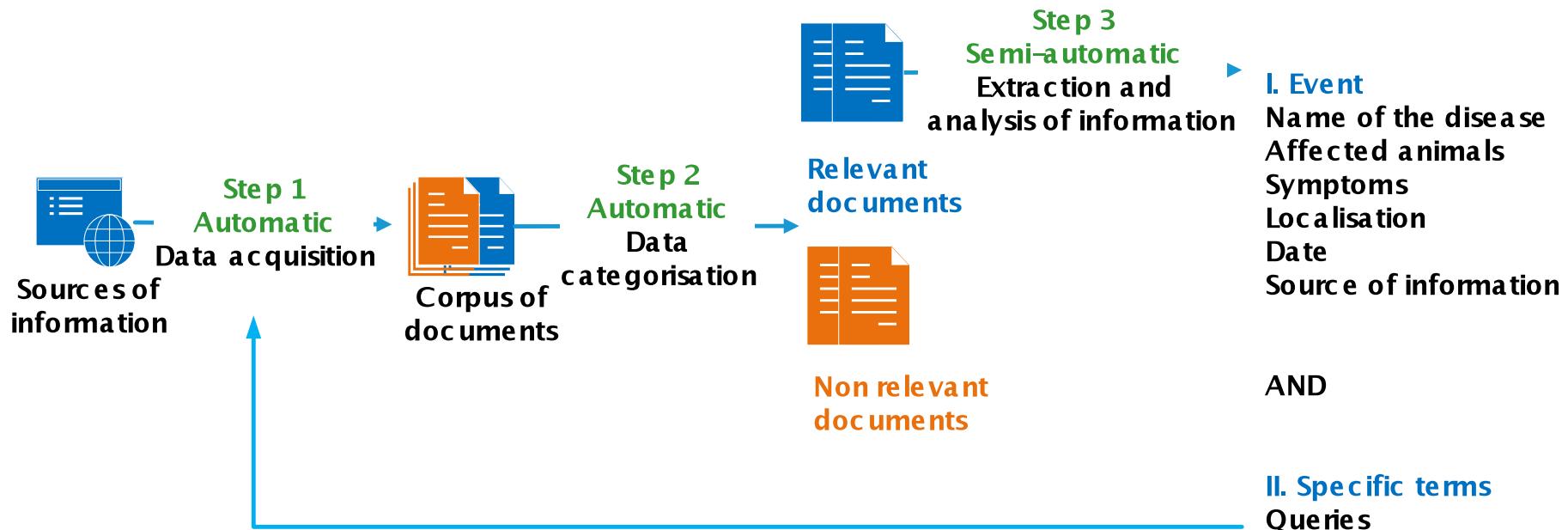
Poland's chief veterinary officer was unavailable for comment, while the county veterinary officer said a statement on the issue will be published later on Wednesday. (Reporting by Anna Wlodarczak-Semczuk; Writing by Marcin Goettig)



# Methodology



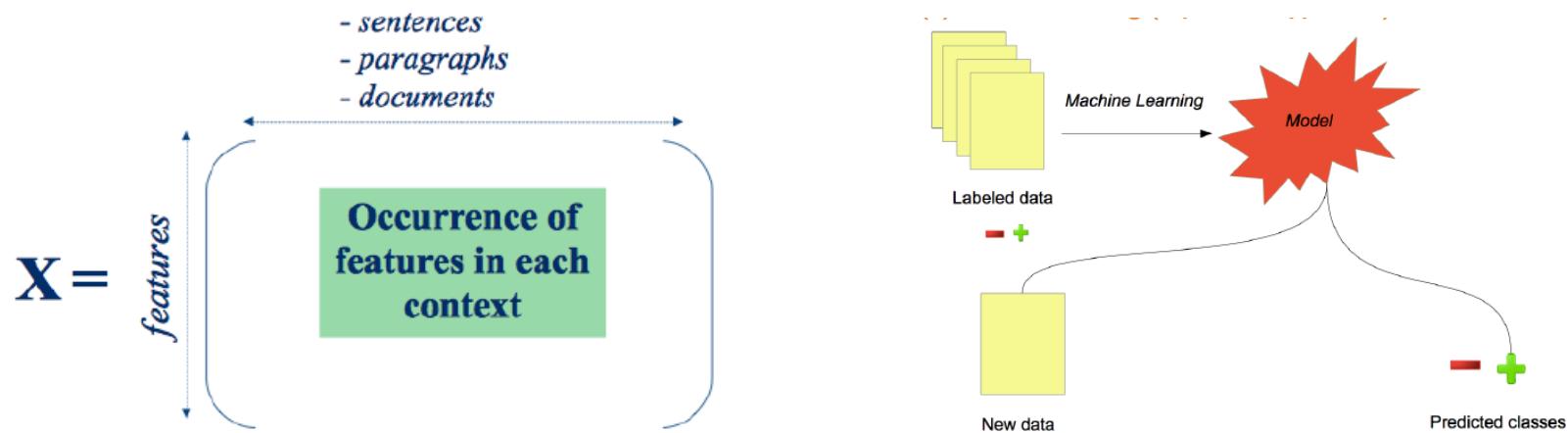
- Step 1: Data acquisition



<https://news.google.com/news/feeds?pz=1&cf=all&ned=en&q=Blue+tongue&output=rss>



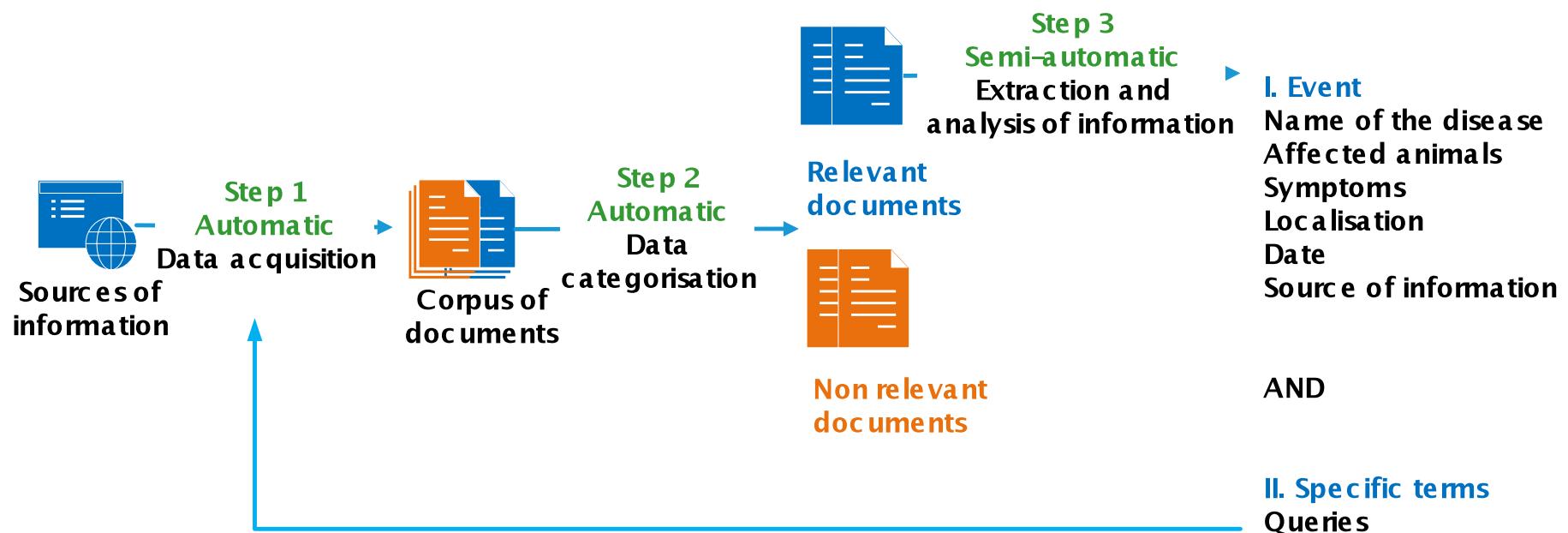
- Step 2: Data classification



Classification algorithm		Naïve Bayes			Support Vector Machine		
Performance		Recall	Precision	F-score	Recall	Precision	F-score
Class	<i>disease</i>	0.724	0.766	0.744	0.657	0.68	0.669
	<i>economy</i>	0.478	0.530	0.503	0.489	0.726	0.584
	<i>general</i>	0.860	0.804	0.831	0.864	0.763	0.810
Weighted average		<b>0.750</b>	<b>0.745</b>	<b>0.747</b>	<b>0.732</b>	<b>0.729</b>	<b>0.725</b>



- Step 3: Information extraction and management**



- **Step 3: Information extraction (I)**

**Aim:** Automatically detecting **key information** from **Web** news articles (country, species, diseases, number of cases, dates, ...)

“Since its initial appearance in **Poland** in **February 2014**, **72** cases of **African Swine Fever** have been detected in **wild boars** and there have been three outbreaks in **pigs**.” - <http://www.thenews.pl>

- Use **dictionaries** (Geonames, HeidelTime, disease names, species names, etc.), and data mining techniques in order to learn extraction rules

**Rules associated with ‘case numbers’:**

(number)(species\_name,1-3) with frequency **26%** and confidence **83%**

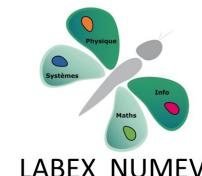
(number)(species\_name,1-2) with frequency **21%** and confidence **100%**



- **Step 3: Information extraction (I)**

- First results for the rule-based approach on the annotated corpus
- Classification based on SVM (features are rules)
- 3 classes: correct, incorrect, partial
- 10-fold cross validation

Type	Accuracy (%)
Locations	70.6
Dates	71.2
Diseases	93.6
Cases	78.1
Species	89.5



**Julien Rabatel,  
LIRMM, Numev, France**



- Step 3: Information extraction (I)

## Veille Sanitaire Internationale

Projet VSI   Accueil   Consultation   Paramétrage

Annotation

Jeux de données annotations

532 article(s) Article précédent Article suivant

**1. Rivers govt. eliminates chickens infected by flu 1 / 532**

Date: 2015-01-20

The Rivers State Government said on Tuesday that it had killed hundreds of fowls infected by the Avian Flu in a privately owned farm in Port Harcourt. The Commissioner for Agriculture, Emma Chinda, said that the farm had been quarantined and decontaminated. He also said no human infection had been recorded. "On January 14, we got a report from a farm that was worrisome. The report we got suggested that the farm may have been infected by the highly pathogenic avian influenza. According to the commissioner, samples of the flu were taken to the Veterinary Research Institute in Vom, Plateau State. The result came out on January 17 and it read positive of highly pathogenic avian influenza. On the basis of that, we had to take necessary steps. Apart from quarantining the farm, we had to depopulate the birds in the farm to stop further spread." Thereafter, we decontaminated the farm. We are containing the situation because officials of government and experts are on ground monitoring the situation", he added. Mr. Chinda said there was no need to panic because government was well equipped to handle the situation. He said before the outbreak, they received information from the Federal Ministry of Agriculture on avian influenza in Kano and a bird market in Lagos. "We were very much on alert and when it happened here, we handled the situation", he said. (NAN)

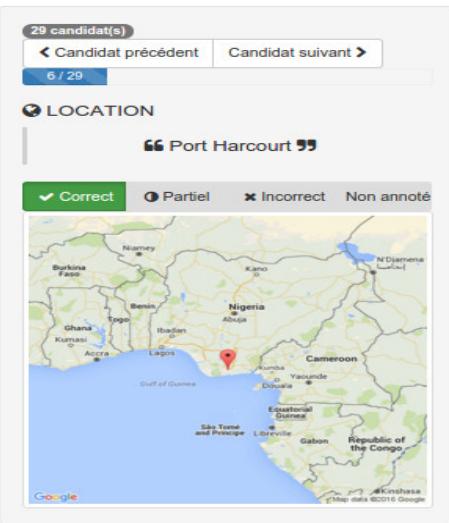
Tout a été annoté Nouveau cas Bilan Non-pertinent

29 candidat(s) Candidat précédent Candidat suivant 6 / 29

LOCATION

Port Harcourt

Correct Partiel Incorrect Non annoté

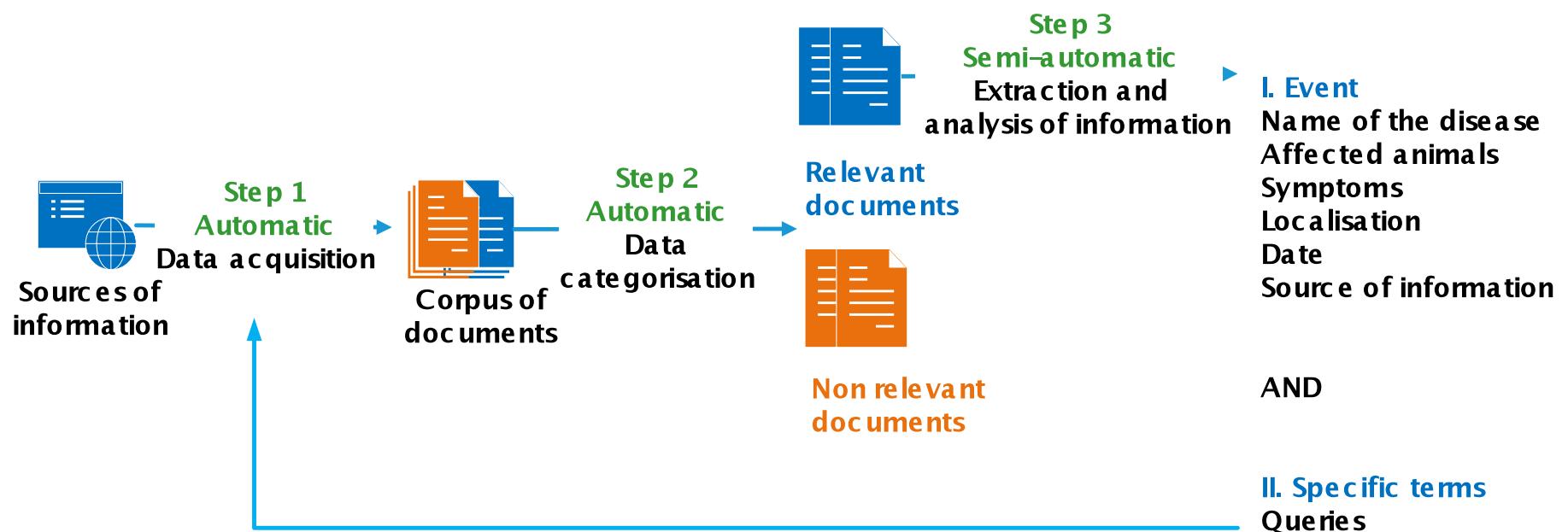


Article précédent Article suivant

Max Devaud-Thomas Fillo Copyright 2015 All Right Reserved.

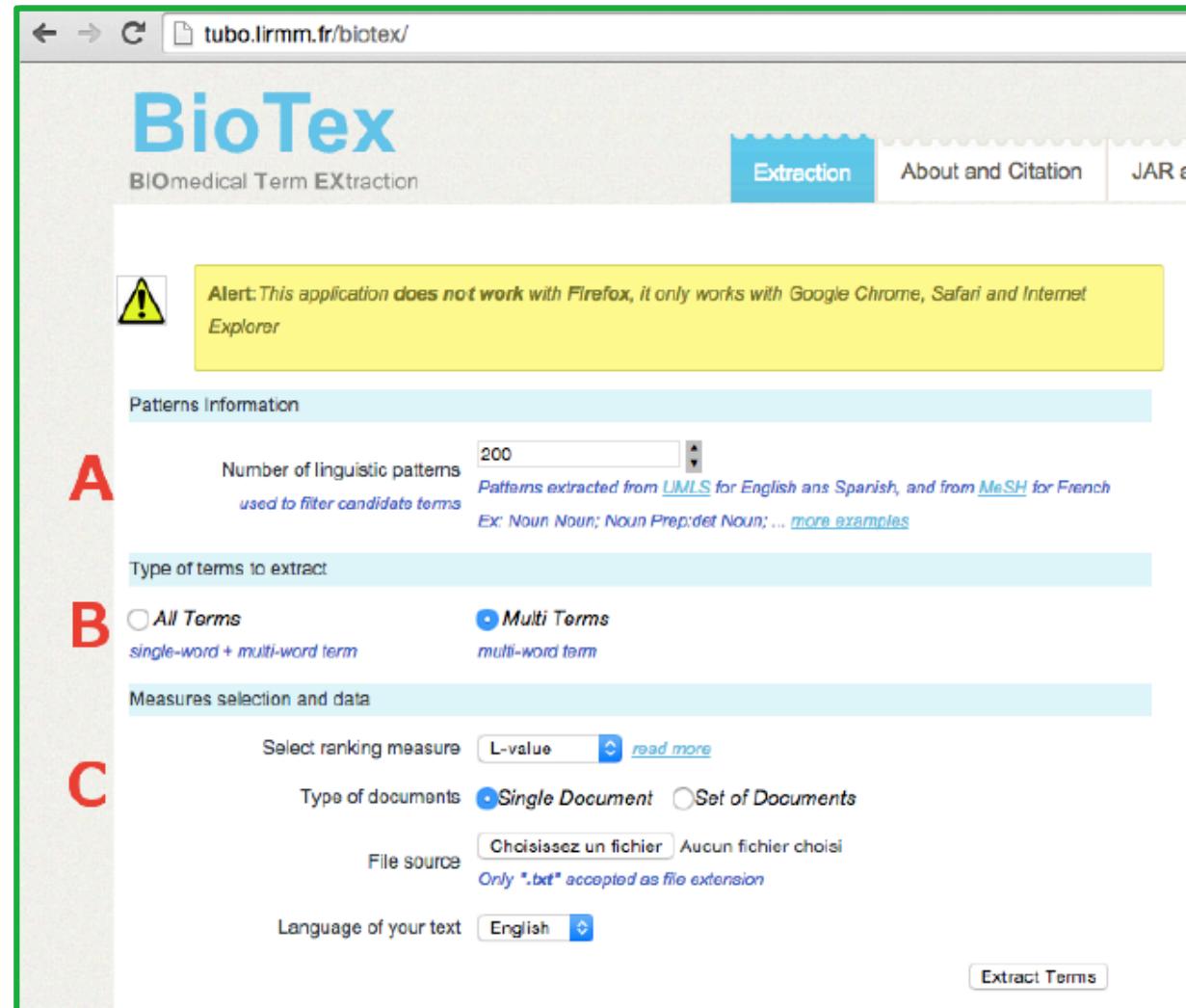


- Step 3: Information management (II)



- II. Querying the Web: (a) *Terminology extraction*

SIFR  
project



The screenshot shows the BioTex web application interface. At the top, there is a warning message: "Alert: This application does not work with Firefox, it only works with Google Chrome, Safari and Internet Explorer". Below this, there are three main configuration sections labeled A, B, and C.

- A Patterns Information:** A dropdown menu shows "200" patterns. Below it, text reads: "Patterns extracted from UMLS for English and Spanish, and from MeSH for French. Ex: Noun Noun; Noun Prep:det Noun; ... [more examples](#)".
- B Type of terms to extract:** A radio button group has "Multi Terms" selected. Below it, text reads: "multi-word term".
- C Measures selection and data:** A dropdown menu for "Select ranking measure" shows "L-value". Below it, text reads: "read more". Another dropdown menu for "Type of documents" has "Single Document" selected. Below that, a file input field says "Choisissez un fichier" and "Aucun fichier choisi". Below that, text reads: "Only \*.txt accepted as file extension". A dropdown menu for "Language of your text" shows "English". At the bottom right is a "Extract Terms" button.



- II. Querying the Web: (b) *Terminology ranking*

### Statistics

- Frequency (TF) → **important** word

$$TF_{i,j} = \frac{n_{i,j}}{\sum_k n_{k,j}}$$

- Inverse Document Frequency (IDF) → **discriminant** word according the distribution in the corpus

$$IDF_i = \log \frac{|D|}{|d_j : t_i \in d_j|}$$

- Global value:

$$TF-IDF_{i,j} = TF_{i,j} \times IDF_i$$



- **II. Querying the Web:** (b) *Terminology ranking*

- BioTex Ranking [Lossio Ventura *et al.* IRJ'2016]:

$$LIDF\text{-}value(t) = P(t_{dom}) \times IDF(t) \times C\text{-}value(t)$$

- A new ranking function to take into account the heterogeneity of the sources ( $S_i$ ) [Arsevska *et al.* CEA'2016]:

$$w(t) = \sum \alpha_i \times \frac{1}{rank_{S_i}(t)}$$

with  $\alpha_i \in [0,1]$  and  $\sum \alpha_i = 1$



- II. Querying the Web: (c) **Terminology validation**

List of Extracted Terms			
See the terms		All	Download XML File
	Number of terms : 1200	Page 1 of 10	1 2 3 4 5 6 7 8 9 10
6	lymph node <small>Validated by : UMLS</small>	<input checked="" type="checkbox"/> YES	<input type="checkbox"/> NO
7	cancer screening <small>Validated by : UMLS</small>	<input checked="" type="checkbox"/> YES	<input type="checkbox"/> NO
8	radiation therapy <small>Validated by : UMLS</small>	<input checked="" type="checkbox"/> YES	<input type="checkbox"/> NO
9	cancer patients	<input checked="" type="checkbox"/> YES	<input checked="" type="checkbox"/> NO
10	endocrine tumors	<input checked="" type="checkbox"/> YES	<input type="checkbox"/> NO
11	cutaneous t-cell lymphomas <small>Validated by : UMLS</small>	<input checked="" type="checkbox"/> YES	<input type="checkbox"/> NO
12	renal cell carcinoma <small>Validated by : UMLS</small>	<input checked="" type="checkbox"/> YES	<input type="checkbox"/> NO

Using of a **Delphi method** [Arsevska et al. LREC'2016]

*Delphi method is to reach group consensus with experts (5 to 7 experts for each disease) when knowledge is not sufficient for a given scientific question*



- II. Querying the Web: (d) ***Association of terms***

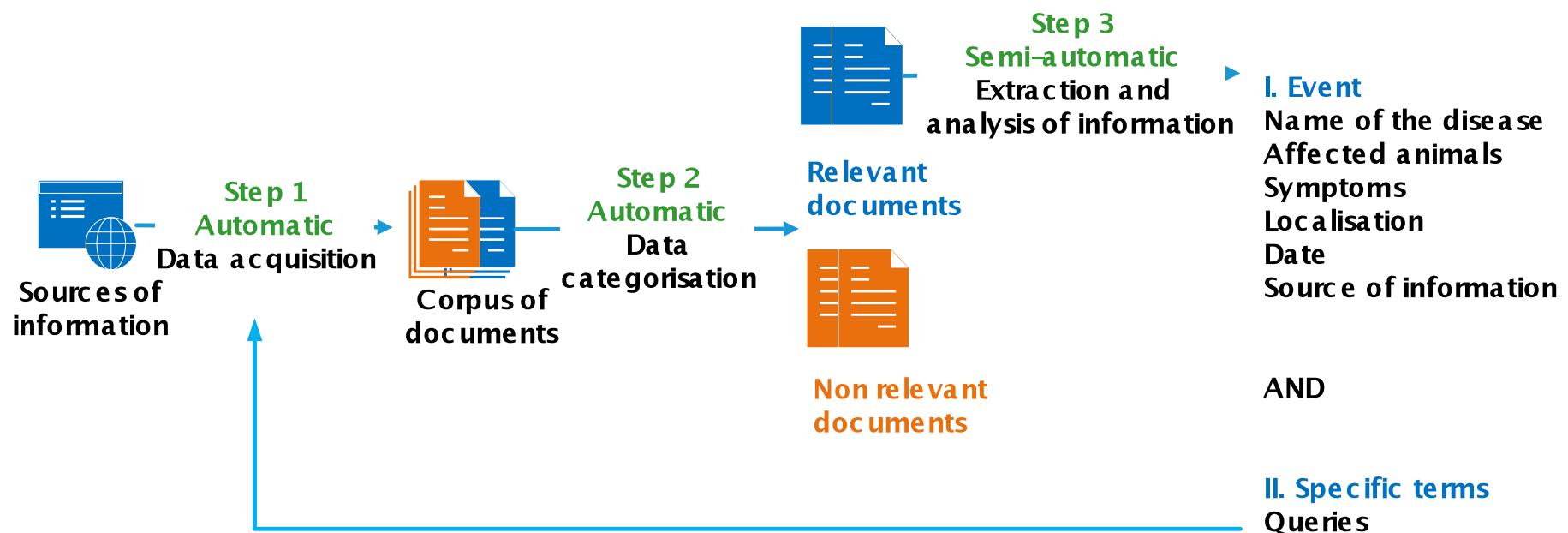
$$D^{AND_{web}} = \frac{2 \times hit(h \text{ AND } cs)}{hit(h) + hit(cs)}$$

[Roche and Prince Informatica'2010 ; Arsevska et al. IJAEIS'2016]

Rank	<b>Bluetongue</b> <i>hosts / clinical signs</i>	<b>Schmallenberg virus infection</b> <i>hosts/ clinical signs</i>
1	general clinical signs / pregnant ewes	stillborn bovine foetuses / camels
2	livestock deaths / sheep	stillborn bovine foetuses / bison
3	embryonic death / cow	aborted foetuses / sheep
4	general clinical signs / sheep	deformed offspring / sheep
5	livestock deaths / cow	stillborn bovine foetuses / deer
6	livestock deaths / deer	aborted foetuses / cattle
7	fever outbreak / sheep	deformed offspring / cattle
8	embryonic death / sheep	stillborn bovine foetuses / calves
9	fever outbreak / cow	deformed offspring / lambs
10	embryonic death / pregnant ewes	acute bronchopneumonia / bison



# Methodology



## Scientific results in Epidemiology:

- **Spatial detection** larger than official surveillance (e.g. PPA detected in Uganda in February and March 2016)
- Detection of **important information** related to diseases not given by official surveillance systems (e.g. preventive vaccinations)
- Detection of information about **other diseases** (e.g. SBV)

## Scientific results in Computer Science:

- New **visualisation** systems
- New **ranking measures** to extract keywords from texts
- New **methods to extract information** in textual documents (i.e. new features based on sequential patterns)
- Implemented **system** in collaboration between LIRMM, TETIS, and CMAEE





**Part 3**

## Applications in agricultural domain

### Sentiment analysis



## Methods in order to identify sentiments: *Towards a sentiment lexicon*

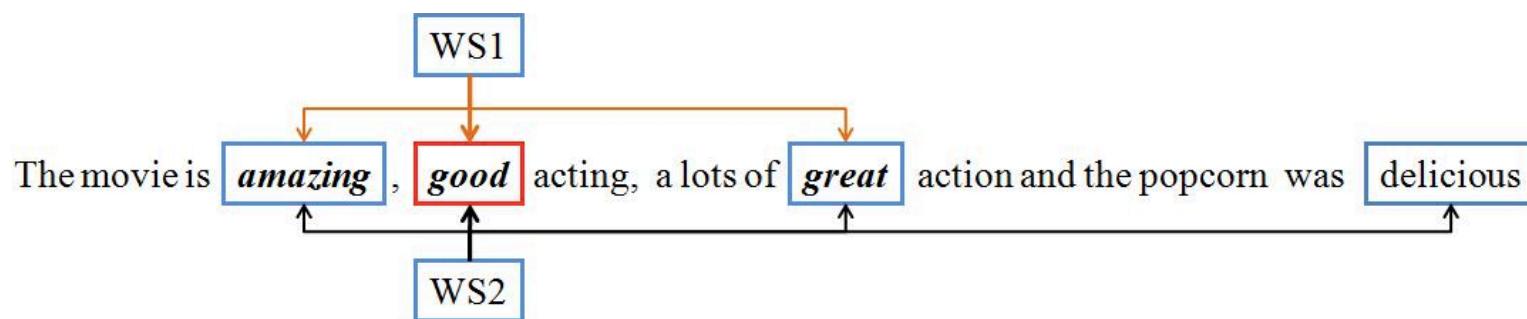
### Step 1: choice of seeds related to opinions

$P = \{good; nice; excellent; positive; fortunate; correct; superior\}$

$N = \{bad; nasty; poor ; negative; unfortunate; wrong; inferior\}$

Construction of 14 corpora related to a specific domain

### Step 2: PoS, Association rules, choice of a window



## Step 3: Statistic selection and web mining

Statistic measures that consist of measuring the association between **seed adjectives** and **candidate adjectives** association based on "hits" from the web (i.e. search engine) and contextual information

**Examples of learnt adjectives:**

*great, hilarious, funny, happy, perfect, important, beautiful, amazing, complete, major, helpful*

→ **Agriculture domain:**

*{gmo; agricultural biotechnology; biotechnology for agriculture}*



**Examples of learnt adjectives:** *green, healthy, enthusiastic, creative, etc.*



**Laura Vanessa Cruz, San Agustin University, Peru**

## Work in progress: tweets and SMS (88milSMS corpus)



Corpus « 88milSMS »

© 2014 Panckhurst, Détrie, Lopez, Moïse, Roche, Verine



Présentation Accès au corpus Références & liens Contact



Présentation



Une équipe pluridisciplinaire de linguistes et d'informaticiens (Rachel Panckhurst, Catherine Détrie, Cédric Lopez, Claudine Moïse, Mathieu Roche, Bertrand Verine (Praxiling, Lirmm, Lidilem, Tétis, Viseo) a recueilli **plus de 88 000 SMS authentiques en français** à Montpellier, en 2011. Cette collecte a été effectuée dans le cadre du projet sud4science LR (Sud4science Languedoc Roussillon. Mutation des pratiques scripturales en communication électronique médiée (financement principal : MSH-M)), lui-même faisant partie du projet international sms4science, coordonné par le CENTAL à l'Université catholique de Louvain (UCL) en Belgique. Lors du recueil des SMS, un questionnaire sociolinguistique a également été proposé aux participants. Les SMS du projet sud4science LR ont été ensuite anonymisés de manière semi-automatique (en collaboration avec des étudiants stagiaires et un juriste-CIL, Nicolas Hvoinsky, SAII, Université Paul-Valéry), puis partiellement transcodés (en français standardisé) et annotés (*cf. Panckhurst et al. 2013*).

- Impact of **repetition of characters** for sentiment analysis: « *j'aaaadore IC'2016 !* » [Khiari et al., JADT'2016]
- Identification of new **spatial entities** and **new spatial relations**: « *IC'2016 est organisé sur montpeul !* » [Zenasni et al., TALN'2016]





## Part 3

# Applications in agricultural domain

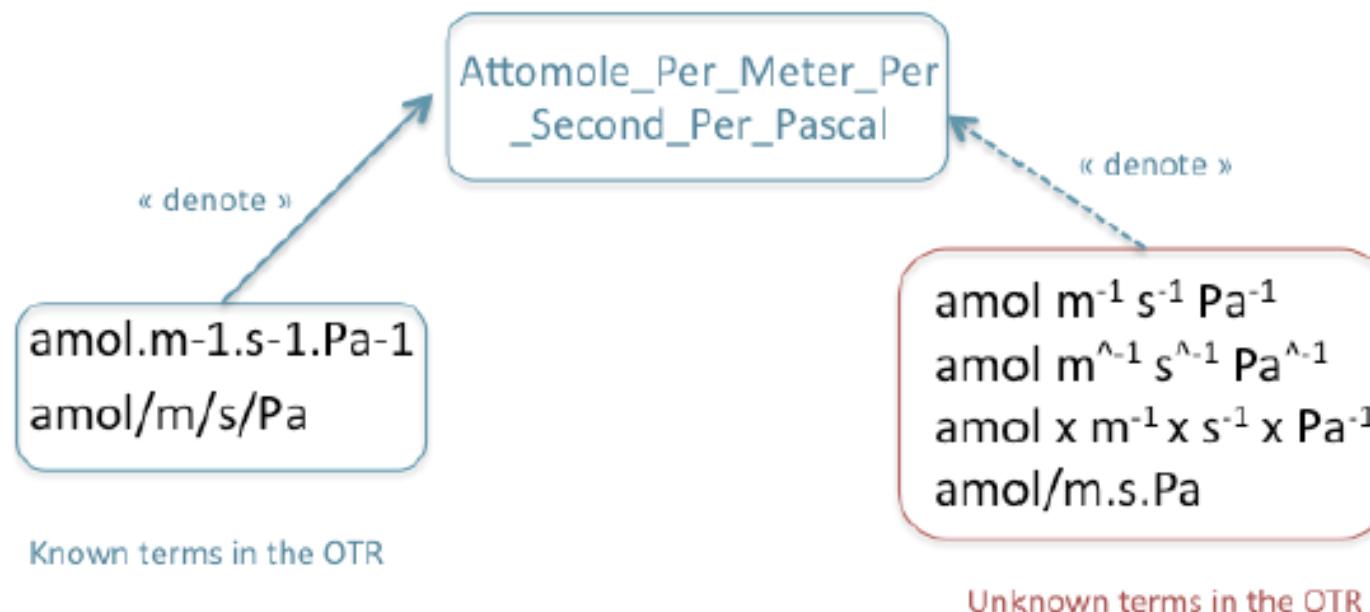
## Information Extraction from experimental data



Aim: **Knowledge management** in food science domain

**Challenging issue: Unit recognition and extraction**

[Berrahou *et al.* KDIR'2013]



## Method:

- Locating unit (*machine learning*)
- Extracting unit (*lexical similarity*)

$$SM_{DL}(u1, u2) = \max[0; \frac{\min(|u1|, |u2|) - DL(u1, u2)}{\min(|u1|, |u2|)}]$$

$$\in [0; 1]$$

Variant term	Reference	SMDc	SMDb
10e10 (cm3.m-1.sec-1.Pa-1)	10e10.cm3.m-1.sec-1.Pa-1	0.87	1
10e-14(cm3/m.s.Pa)	10e-14.cm3/(m.s.Pa)	0.89	1
10e-16cm3.cm/cm.cm2.s.Pa	(10e-16cm3.cm)/(cm2.s.Pa)	0.76	0.8
10e18 (mol.m/Pa.sec.m2)	10e18.mol.m/(Pa.sec.m2)	0.87	1
amol.m-1.s-1.Pa-1	amol.s-1.m-1.Pa-1	0.88	0.75
amol/m.s.Pa	amol/(m.s.Pa)	0.84	1
amol/m.sec.Pa	amol/(m.s.Pa)	0.69	0.75
cm3.um/m2.d.kPa	cm3. $\mu$ m/(m2.d.kPa)	0.77	0.8



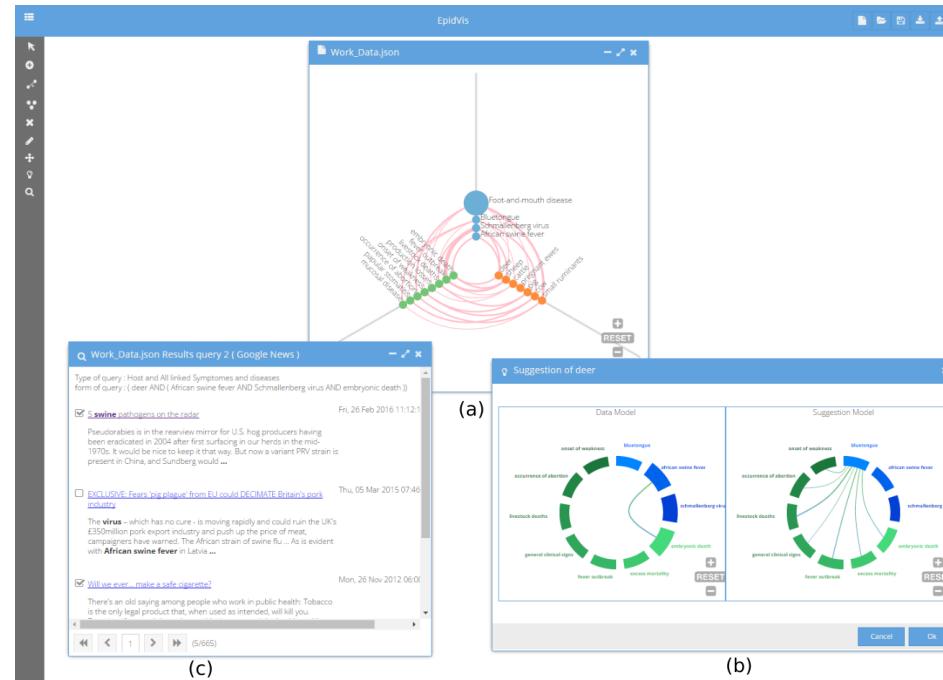
## Part 4

# Conclusions and future work



## New challenges of *Texual Data Science*:

- Matching different types of documents (image/text, video/text, and so forth)
- Integration of visual analytics skills [Fadloun, Inforsid'2016]



## New challenges of *Texual Data Science*:

- Towards the valorisation of data: **Open Data, Data papers**
- From **Textual Data Science** to **Data Science** in pluridisciplinary context:

### **Unification of concepts and associated methods**

**For instance:**

**NLP:** *n-grams, skip-grams, multi-word terms, syntactic relations, etc.*

**Data mining:** *association rules, sequential patterns, etc.*

**Linguistics:** *collocations, etc.*



**1) Titre :** Intégration et visualisation de données issues du projet Patrimoine Numérique Scientifique du Cirad

**Encadrants :** Sandrine Auzoux, Sophie Fortuno et Mathieu Roche

**2) Titre :** Titrage automatique des thématiques identifiées dans les corpus

**Encadrants :** Mathieu Roche, Pascal Poncelet, Julien Velcin et Christophe Gravier

**3) Titre :** Détermination des itinéraires migratoires contextualisés à partir de récits de vies

**Encadrants :** Mathieu Roche, Maguelonne Teisseire et Nelly Robin

+ 2 stages Pro liés aux travaux en épidémiologie animale

