

Exploitation de la spatialité sur des données textuelles hétérogènes provenant de Madagascar

Journée de Variété des Données SHS
28 mai 2018

Jacques Fize, *Doctorant*, TETIS

Dirigé par Mathieu Roche et Maguelonne Teisseire



Sommaire

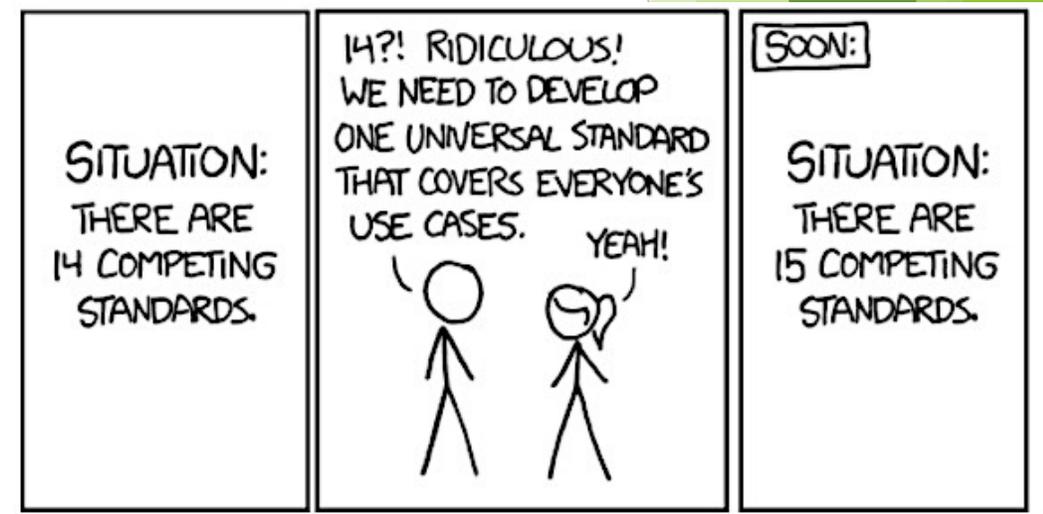
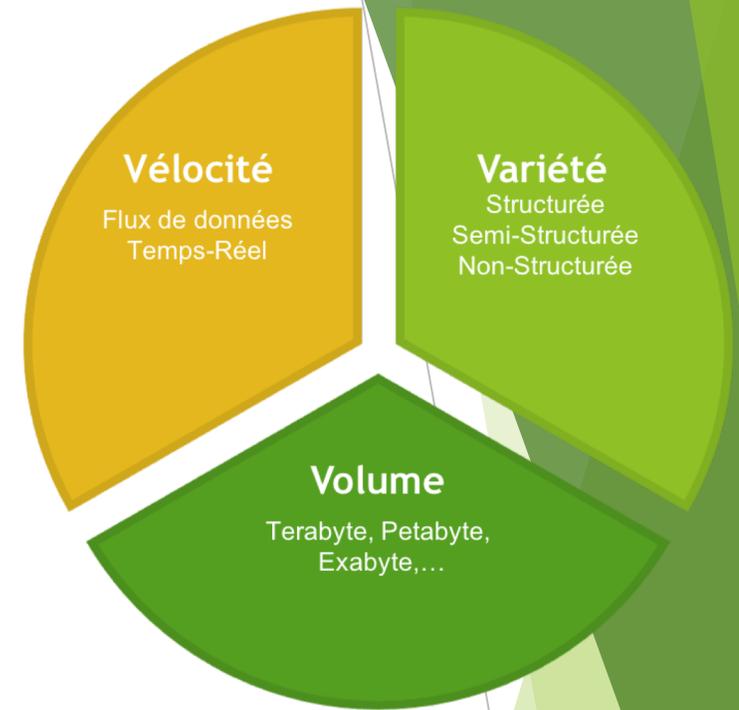
1. Introduction
2. Spatialité dans les textes
3. Comparaison spatiale
4. Données hétérogènes
5. Expérimentations
6. Conclusion

Introduction

The background features a complex, abstract design of overlapping, semi-transparent green triangles and polygons. The colors range from light, pale greens to deep, dark forest greens. The shapes are layered, creating a sense of depth and movement. The overall composition is modern and clean, typical of a corporate or academic presentation slide.

L'ère du Big data

- ▶ Caractérisé par 3 dimensions
 - ▶ Volumétrie, **Vélocité** et **Variété**
 - ▶ ... mais aussi : **Véracité**, **Variabilité**
- ▶ Augmentation constante de l'hétérogénéité de la data
 - ▶ Modernisation des différents secteurs
 - ▶ Évolution des formats en fonction des besoins
 - ▶ Unification des formats ...



Exploiter l'hétérogénéité

- ▶ Augmenter le volume de données exploitables
 - ▶ Recherche d'Information
 - ▶ Extraction d'Information
- ▶ Différentes représentations d'une même information
 - ▶ Synthèse automatique
 - ▶ Découverte de connaissance

Question de recherche

► Quelles méthodes proposer pour mettre en correspondance des données hétérogènes et textuelles ? Approches exploitant différentes dimensions :

- Thématique
- Spatialité
- Temporalité

A	B	C	D	E	F	G	H
Période	Production d	Production d	Production e	Production en %	par rapport au total (washed)		
2000-2001	10 212,48	8 286,06	55,21	44,79			
2001-2002	11 127,00	5 068,14	68,71	31,29			
2002-2003	28 235,75	7 990,08	77,94	22,06			
2003-2004	5 083,73	5 89,62	89,61	10,39			
2004-2005	32 246,04	5 956,92	84,41	15,59			
2005-2006	5 234,40	932,22	84,88	15,12			
2006-2007	22 433,40	7 512,90	74,91	25,09			
2007-2008	4 949,16	2 998,68	62,27	37,73			

REVIEW ARTICLE

Quality assessment of Arabica and Robusta green and roasted coffees – A review

Natalina Cavaco Bicho¹, Fernando Cebola Lidon^{2*}, José Cochicho Ramalho¹ and António Eduardo Leitão¹

¹Grupo Interações Planta-Ambiente (PlantStress), Centro Ambiente Agricultura e Desenvolvimento (BioTrop), Instituto de Investigação Científica Tropical, I.P., Avenida da República, Quinta do Marquês, 2784-505 Oeiras, Portugal

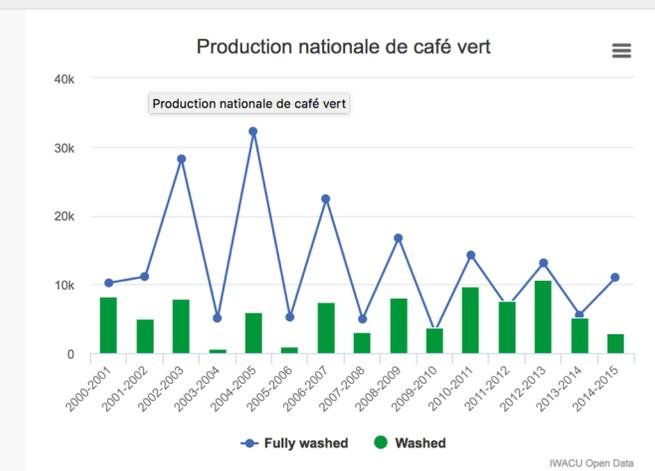
²Departamento de Ciências da Terra, Faculdade de Ciências e Tecnologia, Universidade Nova de Lisboa, Campus da Caparica, 2829-516 Caparica, Portugal

Abstract

This review is a synopsis on coffee quality assessment of green and roasted coffee beans of *Coffea arabica* and *Coffea canephora*. The particle size, medium sieve, most frequent sieve, share split, cumulative calibration, trade homogeneity, mass of 1000 beans, apparent density, strange bodies and defects, mass losses on drying, olfactory and visual parameters, chromatic parameters, soluble solids, pH and chemical characterization (chlorogenic acids, caffeine and trigonelline) is described and evaluated, considering the most important factor associated to the coffee trade, according to a technological perspective.



Evolution de la production de café au Burundi de 2000 à 2015. Les valeurs s



IWACU Open Data

Spatialité dans les textes

The background features abstract, overlapping geometric shapes in various shades of green, ranging from light lime to dark forest green. These shapes are primarily located on the right side of the slide, creating a modern, layered effect. The text is positioned on the left side of the slide, set against a plain white background.

Quelle spatialité ?

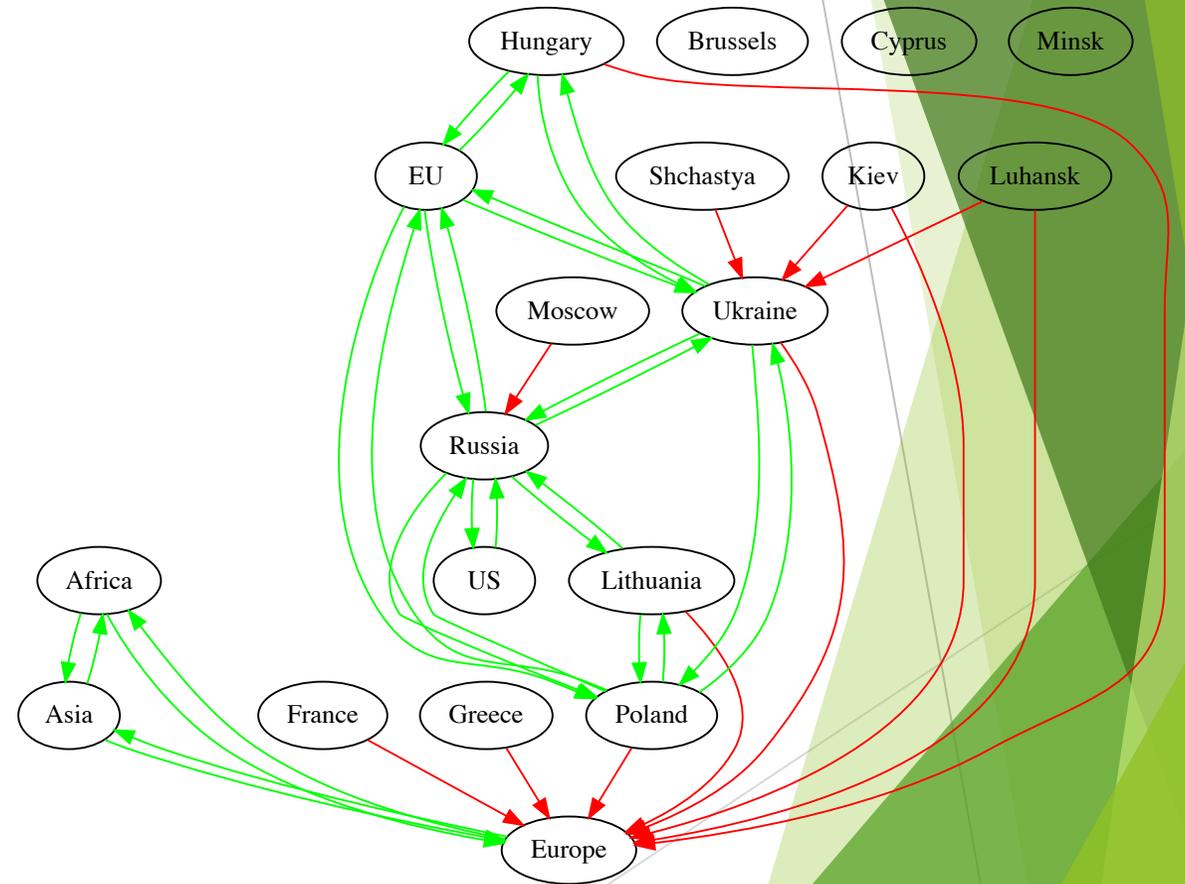
- ▶ La spatialité dans les textes
 - ▶ Entités spatiales : entités localisées dans l'espace
 - ▶ **Relations spatiales** : inclusions, voisinages, ...

```
{
  "en": "Équeurdreville-Hainneville",
  "fr": "Équeurdreville-Hainneville",
  "coord": {
    "lat": 49.648333333333,
    "lon": -1.6547222222222
  },
  "id": "GD4425229",
  "class": [
    "A-ADM4H",
    "A-ADM4"
  ]
}
```

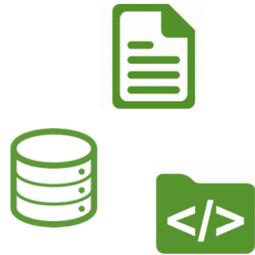


STR ou Spatial Textual Representation

- ▶ Une structure de type **graphe** représentant la **configuration spatiale** d'un document
- ▶ **Entités spatiales** → *nommées* et *référéncées* dans un jeu de données géographiques ou *gazetier* (Geonames, OSM, Dbpedia, ...)
- ▶ **Relations spatiales** → **inclusion** et **adjacence**
- ▶ Perspective
 - ▶ Intégrer des relations thématiques entre les entités spatiales



Construction d'une STR



Document

Extraction de
texte



France, 2018. La ville de Paris reçoit le président de Madagascar. [...] Le maire de Caen, originaire de la ville de Rots était invité.

Identification
des toponymes



France, 2018. La ville de Paris reçoit le président de Madagascar. [...] Le maire de Caen, originaire de la ville de Rots était invité.

Résolution des
toponymes



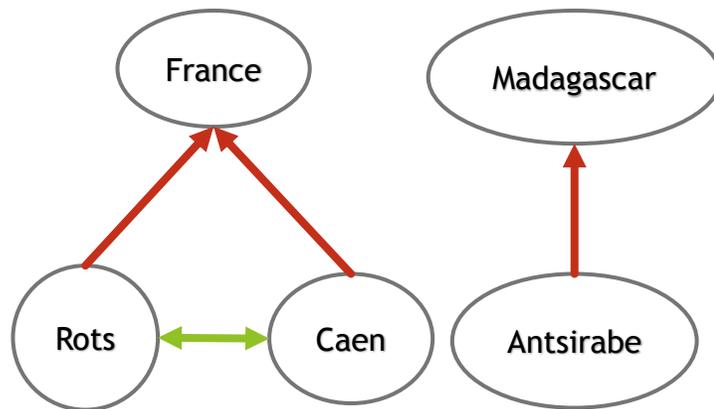
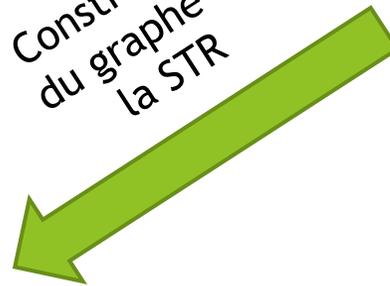
GeoNames



2988507 Paris
1062947 Madagascar
1069166 Antsirabe
...

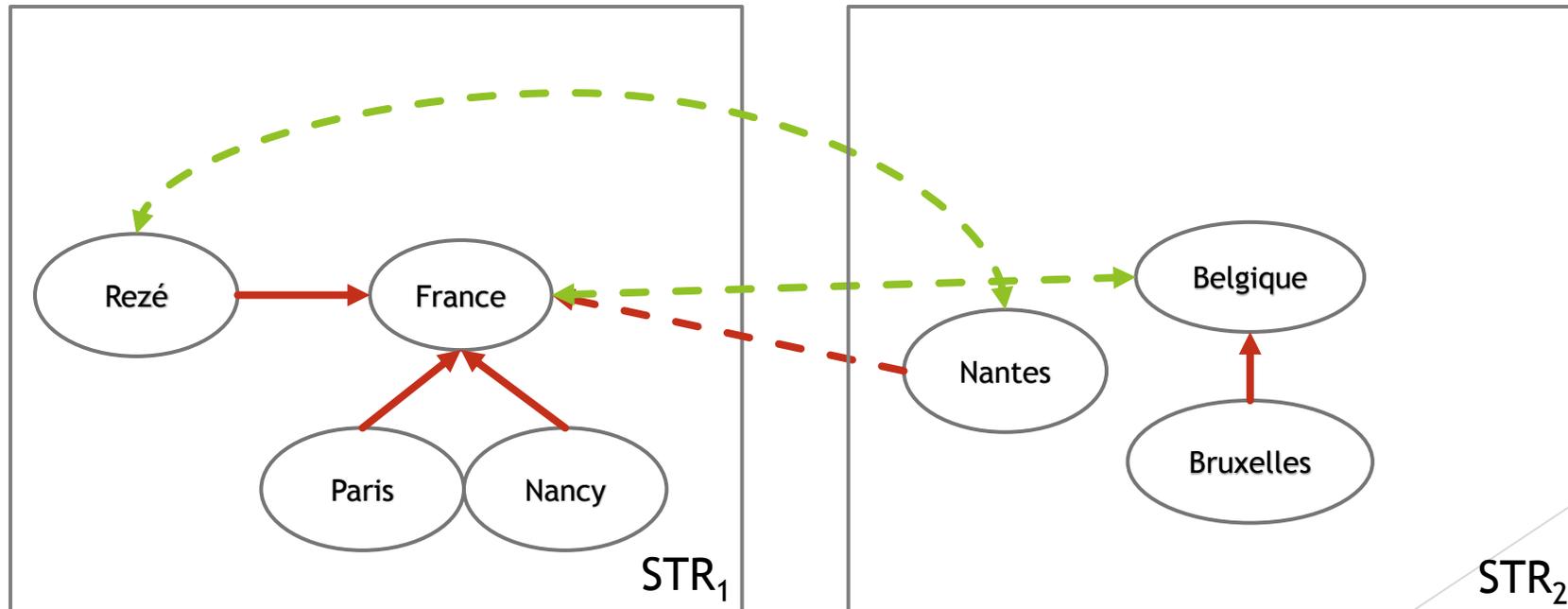


Construction
du graphe de
la STR



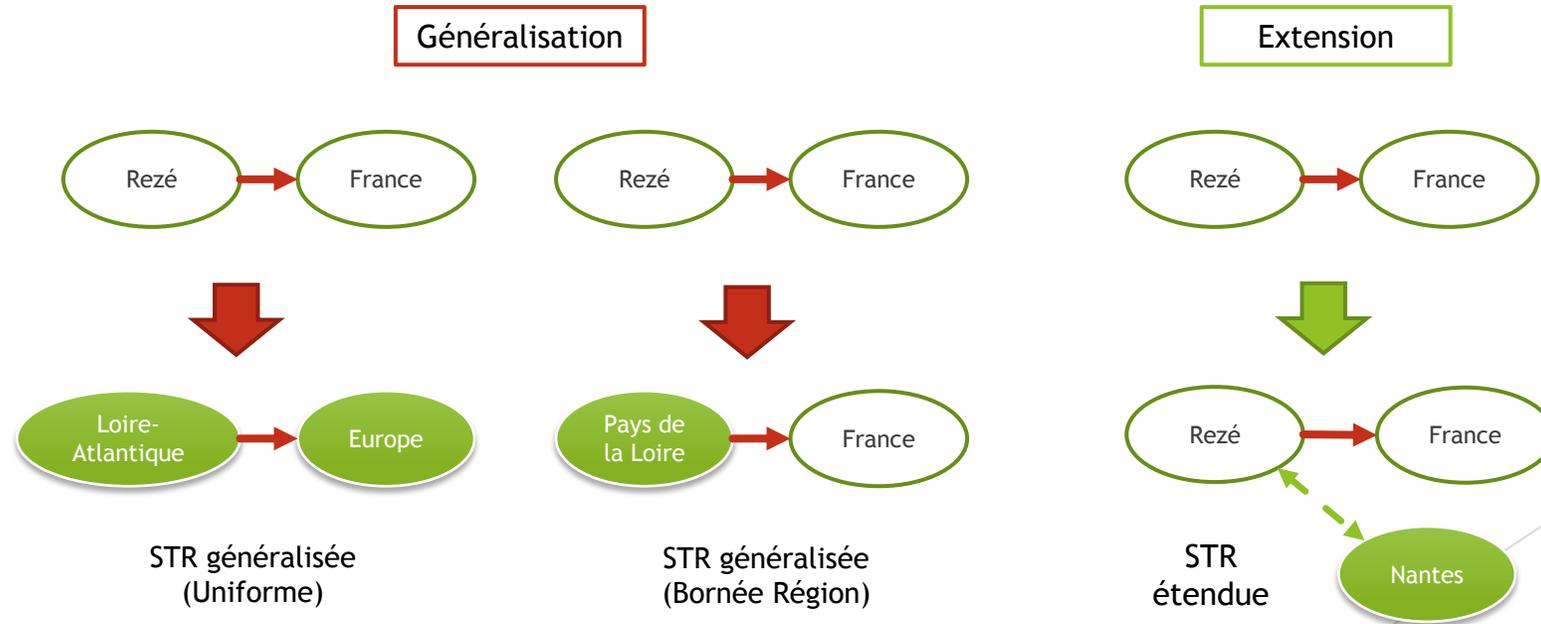
Une information complète ?

- ▶ L'information spatiale peut-être incomplète :
 - ▶ Référentiel incomplet
 - ▶ Élément omis par l'auteur du document (information implicite)
 - ▶ Élément manqué durant le processus de création



Différents types de STR

- ▶ **Objectif** : Augmenter la probabilité de correspondance entre deux STRs
- ▶ Deux transformations → Deux types de STR
 - ▶ Généralisation (uniforme et bornée) → STR Généralisée
 - ▶ Extension → STR Étendue

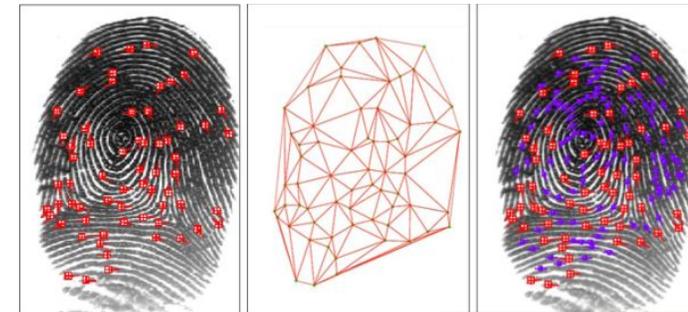
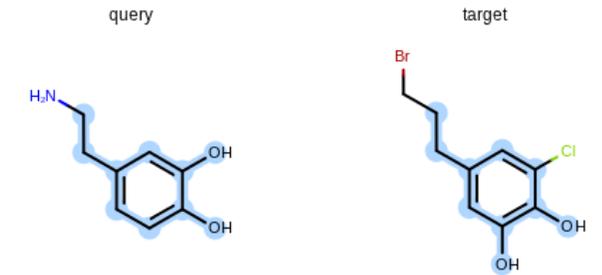


Comparaison spatiale

The slide features a white background with the text 'Comparaison spatiale' in a green, sans-serif font. On the right side, there is a decorative graphic consisting of several overlapping, semi-transparent green triangles and polygons in various shades of green, creating a dynamic, abstract shape.

Comparaison de la spatialité (STR)

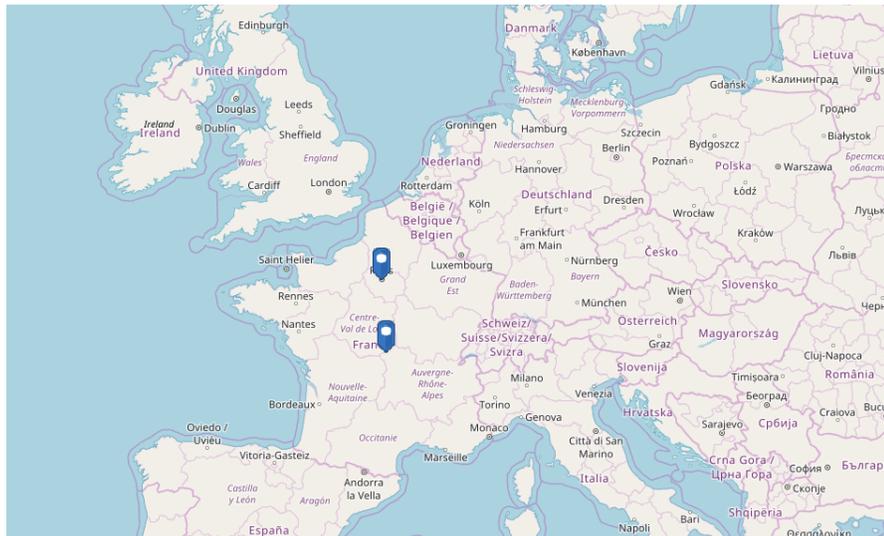
- ▶ Utilisation d'algorithmes de *Graph Matching*
 - ▶ Mesure de similarité entre graphe
 - ▶ Différents algorithmes : MCS, VEO, Graph Edit Distance, Graph Kernels, ...
 - ▶ Implémentation de plusieurs algorithmes dans une librairie Python
- ▶ Comment valider la correspondance entre deux configurations ?
 - ▶ Choix binaire → **Risque de biais**
 - ▶ Solution : **Utilisation d'un ensemble de critères**
 - ▶ Entités partagées ? Entités proches ?
 - ▶ Dispersion des entités similaires ?



Entités Spatiales Communes (ESC) et Entités Spatiales Proches (ESP)

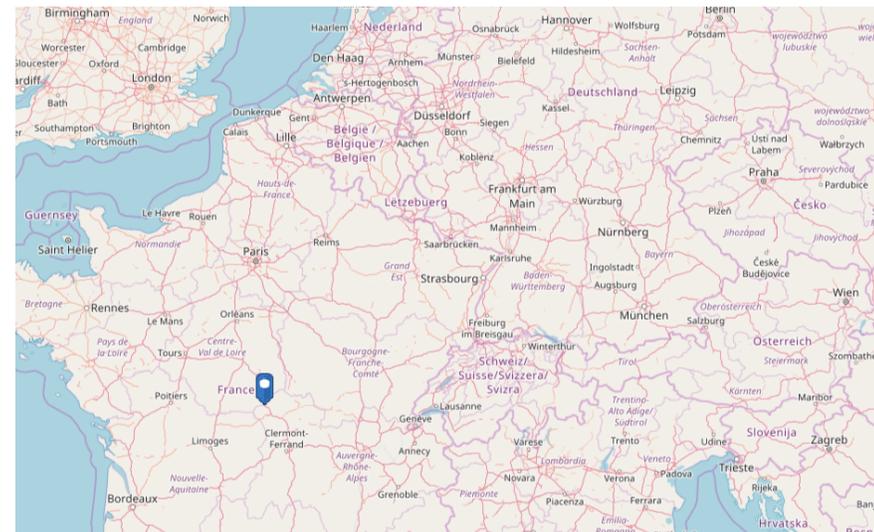
► Entités spatiales communes

- Le critère est validé si les deux STR ont des entités spatiales communes.



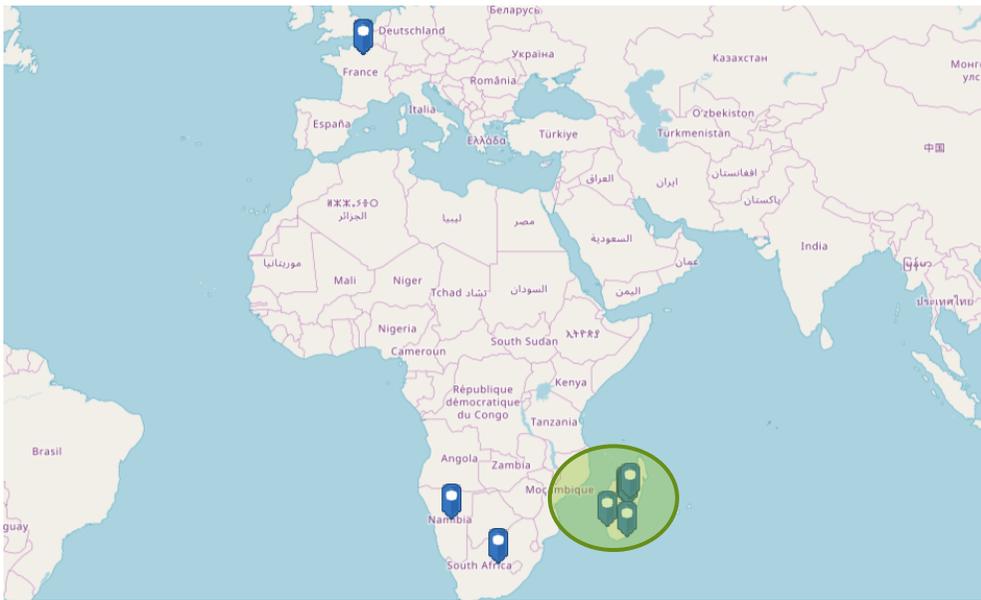
► Entités spatiales proches

- Le critère est validé si un couple d'entités spatiales, dans deux STR sont proches.
- Proximité : relation spatiale ou faible distance géodésique



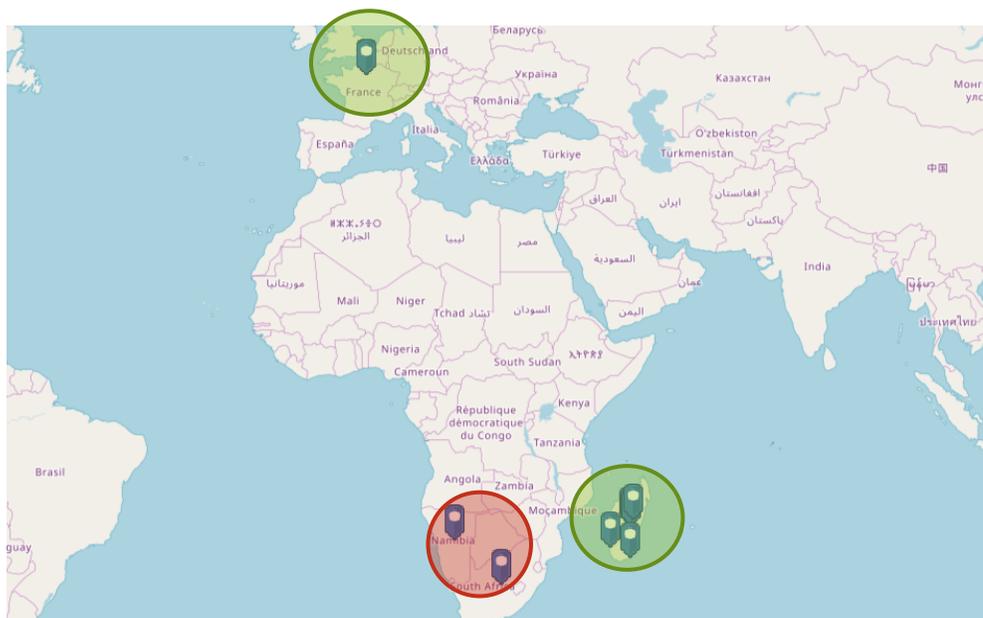
Emprise Spatiale Caractéristique (ESSC)

- ▶ Présence de **groupes significatifs d'entités spatiales proches**.
- ▶ Le critère est validé si : **Un ou plusieurs** de ces **groupes** appartiennent aux **deux STRs**.



L'Emprise Spatiale Stricte (ESS)

- Ce critère est validé si les deux **STRs** partagent une **répartition spatiale identique** de leurs **entités spatiales**.



Données hétérogènes

The right side of the slide features a decorative graphic composed of several overlapping, semi-transparent green triangles and polygons. The colors range from a light, pale green to a dark, forest green. The shapes are arranged in a way that creates a sense of depth and movement, with some shapes appearing to be in front of others. The overall effect is a modern, abstract design element.

Corpus homogène - PADI-Web

- ▶ Corpus¹ construit et annoté durant la thèse de E. Arsevska
 - ▶ **Élaboration d'une méthode semi-automatique pour l'identification et le traitement des signaux d'émergence pour la veille internationale sur les maladies animales infectieuses.** Santé publique et épidémiologie. Université Paris-Saclay, 2017.
 - ▶ Pourquoi utiliser ces données ?
 - ▶ Beaucoup d'informations spatiales
 - ▶ Documents annotés

1. Rabatel Julien; Arsevska, Elena; de Goër de Hervé, Jocelyn; Falala, Sylvain; Lancelot, Renaud; Roche, Mathieu, 2017, "PADI-web corpus: news manually labeled", [doi:10.18167/DVN1/KMTIFG](https://doi.org/10.18167/DVN1/KMTIFG), CIRAD Dataverse, V2

Corpus hétérogènes - AgroMada

- ▶ Ensemble de données produites par le CIRAD
 - ▶ Données fortement hétérogènes tant sur le :
 - ▶ Contenu (langue(s), structure du document, richesse, ...)
 - ▶ Format (Voir Tableau)
 - ▶ Données massives : 13 742 documents (~13 Go)

format	Nombre de doc.
doc	6651
docx	838
html	3465
pdf	1931
ppt	228
pptx	40
sql	1
txt	157
xls	2544
xlsx	126
xml	29



Conséquence de l'hétérogénéité dans le processus de création de la STR

- ▶ Outils de reconnaissance d'entités nommées (NER) : StanfordCoreNLP, NLTK, OpenNLP, Spacy ...
 - ▶ Fonctionne sur des textes homogènes
 - ▶ ... beaucoup moins sur des sources hétérogènes
- ▶ Cause
 - ▶ Méthodes basées sur la syntaxe
 - ▶ Utilisation de règles
 - ▶ Apprentissage → Modèles créés à partir de textes classiques
- ▶ Conséquences
 - ▶ Toponymes non-identifiés
 - ▶ ... ou beaucoup de faux-positifs

Corpus Homogène (PadiWeb)			
NER	P	R	F1
StanfordNER	0.59	0.77	0.67
Polyglot	0.53	0.72	0.61
NLTK	0.42	0.66	0.52
Spacy	0.40	0.65	0.50

Corpus Hétérogène (Agromada)			
NER	P	R	F
StanfordNER	0.31	0.16	0.22
Polyglot	0.20	0.35	0.26
NLTK	0.13	0.15	0.14
Spacy	0.14	0.84	0.25

Exemple



Expérimentations

The background features a series of overlapping, semi-transparent green triangles and polygons of various shades, ranging from light lime green to dark forest green. These shapes are primarily concentrated on the right side of the frame, creating a dynamic, layered effect. The left side of the image is mostly white, providing a clean space for the text.

Expérimentation

- ▶ Objectif : Définir **quelles mesures de similarité** et **types de STR** donnent les meilleures correspondances sur la dimension spatiale.
- ▶ Deux corpus utilisés
 - ▶ Homogène : brèves traitant d'épidémies animales variées
 - ▶ Hétérogène: documents produits à Madagascar (ppt, xls, doc, pdf) traitant de thématiques de recherche variées liées au projet BVLAC¹
- ▶ Mesure de performance
 - ▶ **Précision@n**: pourcentage des STR n-plus similaires qui valide le critère (pour une STR)

1. Projet de Mise en Valeur et de Protection du Bassin Versant du Lac Alaotra)

PADI-WEB - Résultats

- ▶ Mesures respectant le mieux les critères, selon chaque type de STR (**Précision@3**)

Type	ESC	ESP	ESSC	ESS
<i>Normal</i>	BOC	MCS	MCS	BOWSE
<i>Gen_country</i>	VEO	VEO	VEO	VEO
<i>Gen_region</i>	VEO	MCS	MCS	MCS
<i>Extension 1</i>	MCS	MCS	MCS	VEO

- ▶ Scores associés aux meilleurs mesures indiquées précédemment (**Précision@3**)

Type	ESC	ESP	ESSC	ESS
<i>Normal</i>	0,872	0,772	0,424	0,376
<i>Gen_country</i>	0,756	0,872	0,48	0,368
<i>Gen_region</i>	0,904	0,796	0,464	0,396
<i>Extension 1</i>	0,896	0,796	0,456	0,384

- ▶ MCS, VEO, BOC donne les meilleurs correspondances entre STR
- ▶ Des mesures plus complexes donnent des résultats plus faibles
 - ▶ Nature « éparses » des graphes
- ▶ Les transformations de STR améliore de manière significative la qualité des correspondances

Données Madagascar - Résultats

- ▶ Mesures respectant le mieux les critères, selon chaque type de STR (**Précision@3**)

Type	ESC	ESP	ESSC	ESS
<i>Normal</i>	BOWSE/VEO	MCS	BOC	BOC
<i>Gen_country</i>	MCS	MCS	MCS	BOC
<i>Gen_region</i>	MCS/BOWS E/VEO	MCS	BOC	BOC
<i>Extension 1</i>	VEO	VEO	BOC	BOC

- ▶ Scores associés aux meilleurs mesures indiquées précédemment (**Précision@3**)

Type	ESC	ESP	ESSC	ESS
<i>Normal</i>	0,95	0,98	0,71	0,35
<i>Gen_country</i>	0,93	0,95	0,64	0,34
<i>Gen_region</i>	0,95	0,97	0,70	0,34
<i>Extension 1</i>	0,95	0,99	0,73	0,31

- ▶ MCS, VEO, BOC donne les meilleurs correspondances entre STR
- ▶ La transformation améliore la correspondance entre les STR
 - ▶ Moins significatif

Conclusion

- ▶ Proposition d'une représentation dédiée à la spatialité dans les textes
 - ▶ STR, une structure graphe composée d'entités spatiales liées par leur relation spatiale (adjacence, inclusion)
 - ▶ Utilisable sur des données textuelles hétérogènes
- ▶ Similarité spatiale entre documents
 - ▶ Des résultats encourageants !
- ▶ Perspectives
 - ▶ Ajout de la thématique dans la STR

Merci pour votre attention !