

Rapport de stage

**Extraction d'entités nommées géographiques et exploitation dans le cadre
du Web des données**

Étudiant : Arbona Arnaud

Formation : Master 2 Informatique Fondamentale et Ingénierie

Parcours : Ingénierie Spécialité Ingénieur en Web et IA

Encadrante de stage : Madame Anne Toulet, Cirad

Tutrice enseignante : Madame Catherine Faron, Polytech Nice-Sophia



Table des matières

Table des matières	2
1. Introduction.....	4
1.1 Contexte et sujet du stage	4
1.2 Plan du rapport.....	5
2. Description du travail proposé.....	6
2.1 Présentation du projet Issa	6
2.2 Descriptif du stage	7
3. Description du travail réalisé	8
3.1 Planning du stage	8
3.2 Tâches réalisées	9
3.2.1 État de l’art des référentiels géographiques	9
3.2.2 Conception et mise en place de la chaîne de traitement	11
3.2.3 Enrichissement des données et visualisation cartographique	21
Conclusion et perspectives.....	24
Bibliographie.....	25
Table des figures.....	27
Annexes	28
Entity-fishing.py	28
Requête.py	30
Mise_en_confinité_GeoNames_rdf.py	33
Requête Sparql.....	33
Résumé	34
Abstract	34

Remerciements

Je tiens à remercier toutes les personnes qui ont participé au bon déroulement de mon stage et qui m'ont aidé lors de la rédaction de ce rapport.

En premier lieu, je tiens à remercier mon encadrante de stage, Madame Anne Toulet pour son accueil et ses conseils au cours de ce stage. Je remercie également Monsieur Franck Michel de l'équipe Wimmics et Monsieur Andon Tchechmedjiev de l'IMT Mines Alès pour leur aide et leur appui.

Ce travail a été réalisé dans le cadre du projet Issa soutenu par le GIS Collex-Persée (<https://www.collexpersee.eu/projet/issa/>)

1. Introduction

1.1 Contexte et sujet du stage

La science ouverte est un mouvement international qui a pour but de rendre accessible la recherche scientifique et les données qu'elle produit. Dans cet objectif, les archives ouvertes, bases de données documentaires accessibles librement et gratuitement sur internet contenant des documents issus de la recherche scientifique, accentuent leurs efforts pour accroître l'accessibilité aux ressources dont elles disposent.

Ce stage s'inscrit dans le cadre du projet Issa (Indexation Sémantique d'une archive scientifique et Services Associés pour la science ouverte), lauréat de l'appel à projet CollEx-Persée 2019/2020.

L'objectif de ce stage est de permettre un accès et une interopérabilité accrus à des publications scientifiques proposées par une archive ouverte en adoptant des techniques d'indexation sémantique adossée à des référentiels terminologiques standards. Le recours aux techniques du Web sémantique et du traitement naturel des langues sera privilégié.

Dans ce travail, nous nous intéresserons en particulier à la question de l'indexation¹ par des mots-clés géographiques : extraction d'entités nommées² géographiques, alignement avec des référentiels sémantiques standards et visualisation cartographique.

Tout au long de ce projet, l'archive ouverte du cirad, Agritrop (<https://agritrop.cirad.fr/>), servira de cas d'usage.

Agritrop est un portail qui propose essentiellement des publications scientifiques, mais aussi un fonds de cartes et de documents anciens. Chaque document est décrit par des métadonnées riches parmi lesquelles se trouvent des descripteurs thématiques et géographiques issus du thésaurus³ Agrovoc⁴.

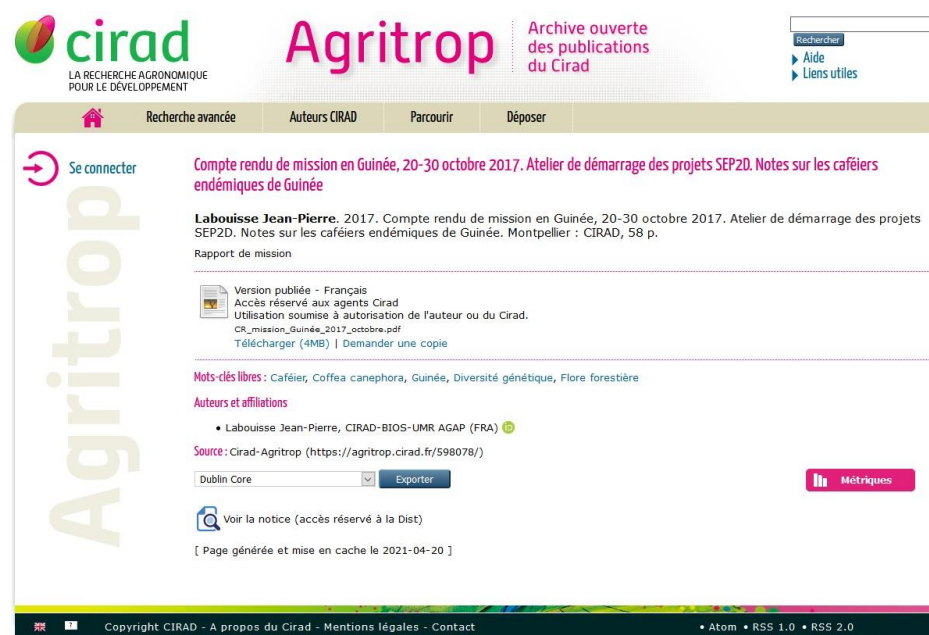


Figure 1 : Exemple d'une notice d'un article dans Agritrop

¹ L'indexation permet de préciser le contenu d'un document à travers des mots-clés et ainsi de retrouver dans un catalogue tous les documents qui traitent d'un sujet donné quel que soit le support.

² Concept d'intérêt dans un document donné

³ Un thésaurus est un lexique structuré de mots-clés permettant l'analyse et le classement de documents

⁴ Thésaurus développé et maintenu par la FAO depuis le début des années 80 (<http://aims.fao.org/fr/agrovoc>)

Le projet ISSA est composé de trois partenaires institutionnels qui travaillent en collaboration :

- Le Cirad (Centre de coopération internationale en recherche agronomique pour le développement), qui est un organisme français de recherche agronomique et de coopération internationale pour le développement durable des régions tropicales et méditerranéennes.
- IMT Mines Alès (Coordinateur Andon Tchechmedjiev).
- L'équipe Wimmics, qui est une équipe de recherche conjointe entre l'INRIA Sophia Antipolis-Méditerranée et I3S – CNRS et Université côte d'Azur – (Coordinateur Franck Michel).

Nous allons maintenant présenter le contenu du rapport.

1.2 Plan du rapport

Dans les chapitres suivants, nous commencerons par une présentation détaillée du sujet de stage.

Puis nous expliquerons le travail réalisé pour atteindre les objectifs explicités dans la partie précédente :

- État de l'art des référentiels géographiques
- Conception et réalisation d'une chaîne de traitement
- Analyse des résultats obtenus
- Enrichissement des données et visualisation cartographique

Dans la conclusion, nous reviendrons sur l'expérience acquise au cours de ce stage, en particulier les domaines scientifiques sur lesquelles j'ai pu me perfectionner ainsi que sur les connaissances et compétences que j'ai pu améliorer.

2. Description du travail proposé

2.1 Présentation du projet Issa

Comme indiqué dans l'introduction, mon stage s'inscrit dans le projet Issa qui se décompose en plusieurs work packages.

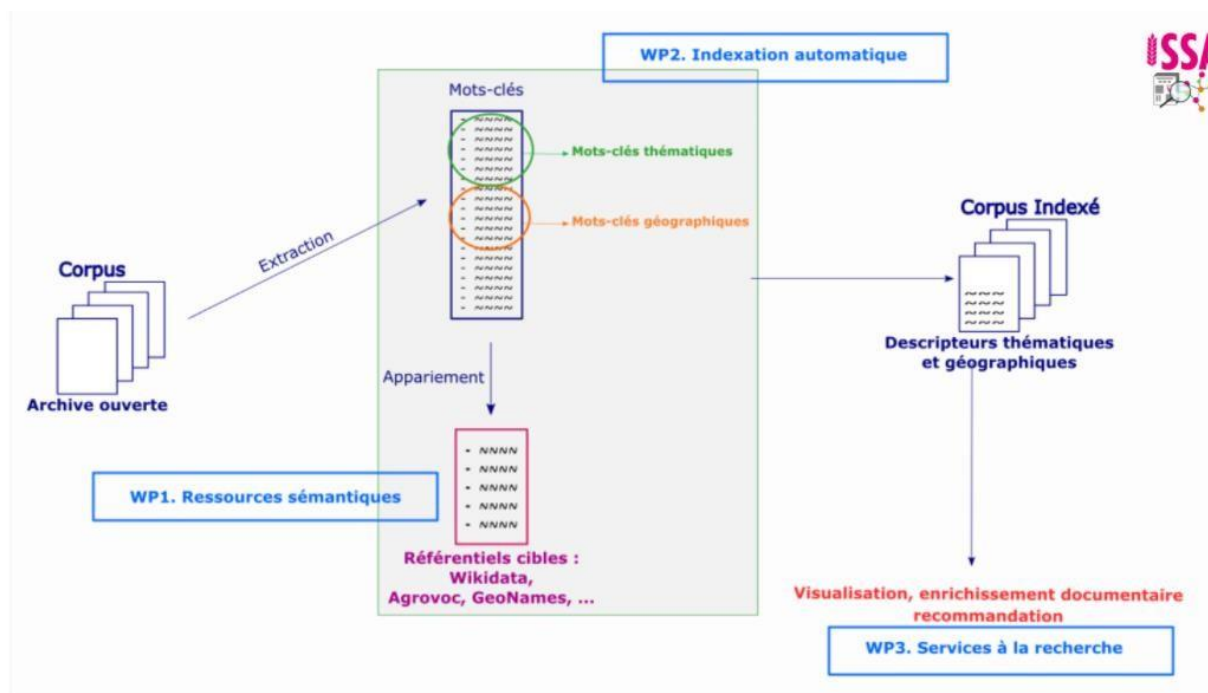


Figure 2 : Décomposition du projet issa

Le WP1 concerne les ressources sémantiques d'intérêt pour le projet dans le domaine de l'agronomie et la géographie. Cette étude sera basée sur les besoins d'indexation dans Agritrop et proposera des pointeurs vers des référentiels terminologiques adaptés à d'autres domaines pour permettre de transférer la démarche à d'autres communautés.

Le WP2 a pour but de concevoir une chaîne de traitement⁵ générique permettant l'indexation automatique d'un corpus scientifique. Il faudra tout d'abord effectuer une extraction structurée complète des textes (traitement d'un pdf pour obtenir un texte structuré exploitable par les machines) puis un liage des entités nommées. Pour cela nous utiliserons des techniques de reconnaissance et désambiguïsation d'entités nommées.

Le WP3 a pour objectif d'exploiter les descripteurs sémantiques en proposant différents services aux utilisateurs, comme une visualisation enrichie. Cela sera possible grâce aux descripteurs globaux ainsi que les entités nommées obtenues à la sortie du WP2.

Dans le paragraphe suivant, je détaillerai plus précisément les tâches à mener dans le cadre de mon stage.

⁵ Série de processus permettant d'extraire de l'information ou produire du savoir à partir de données brutes de manière automatisée

2.2 Descriptif du stage

La problématique de mon sujet de stage se concentre sur l'extraction d'entités nommées géographiques et leur visualisation cartographique.

Actuellement, l'indexation des documents par mots-clés géographiques est effectuée manuellement par des documentalistes.

L'objectif principal est de trouver comment automatiser cette procédure et enrichir les informations avec de la visualisation cartographique.

Pour atteindre cet objectif, plusieurs étapes seront nécessaires : 1) choisir un référentiel géographique pertinent pour le projet ; 2) extraire les entités nommées géographiques des documents et les lier avec le référentiel cible ; 3) proposer un service de visualisation cartographique de ces entités nommées.

Dans ce cadre-là, les tâches attendues sont les suivantes :

1. États de l'art

- Dresser un inventaire détaillé des référentiels géographiques existants en comparant leurs différentes approches. On s'intéressera particulièrement à GeoNames.
- Étudier et comparer les différents outils d'extraction d'entités nommées existants, en particulier pour les entités géographiques.

2. Conception et réalisation d'une chaîne de traitement

- Développer et appliquer une chaîne d'extraction et de désambiguïsation des entités nommées
- Liage des entités nommées extraites avec le référentiel géographique choisi sur la base de l'état de l'art effectué précédemment,

Cette chaîne de traitement sera testée sur un jeu de publications scientifiques issues d'Agritrop.

3. Exploitation de l'indexation sémantique

Exploiter les descripteurs géographiques obtenus en utilisant les technologies du Web sémantique et le Web de données, par exemple en permettant la visualisation géographique (cartographie) et/ou l'enrichissement encyclopédique à partir des mots-clés géographiques obtenus précédemment.

Il sera nécessaire de mettre en application des règles de bonne pratique de développement et d'industrialisation d'une application. Il s'agit de délivrer un produit fini, facilement déployable et configurable dans d'autres environnements.

3. Description du travail réalisé

3.1 Planning du stage

Mois	Tâche
Avril	Étude de la bibliographie + état de l'art
Mai	État de l'art + tests outils d'extraction
Juin	Chaîne de traitement
Juillet	Chaîne de traitement et visualisation cartographique
Août	Documentation + amélioration de la chaîne de traitement et de la visualisation

3.2 Tâches réalisées

3.2.1 État de l’art des référentiels géographiques

Tel qu’exposé précédemment, le premier objectif de mon stage est de réaliser un inventaire des référentiels géographiques.

Tout d’abord, sur la partie concernant l’état de l’art des référentiels géographiques, il a fallu déterminer quels sont les critères de sélection qui permettront de choisir un référentiel adapté aux besoins du projet.

Les principales caractéristiques recherchées sont les suivantes :

- Présence des coordonnées géographiques : latitude, longitude
- Disponibilité au format du web sémantique
- Ressource en libre d’accès
- Accessible via une API Rest / Sparql endpoint
- Couverture géographique mondiale

Ci-dessous vous trouverez un tableau récapitulatif de l’analyse effectuée en fonction de ces critères.

Les liens vers ces référentiels sont indiqués dans la partie « Liens utiles » à la fin de ce document.

Nom	Développeur	Objectif Projet	Dernière date de mise à jour	Licence	API	Type Donnée	Projet utilisant ces référentiels
GeoNames	Wiki + équipe d’aide aux développeme nt	géographique	2021	Libre de partager et de modifier Préciser que l’on utilise GeoNames	oui	Json xml/rdf	Multimap microsoft Popfly
LinkedGeoData OpenStreetMap	AKWS (Agile Knowledge Engineering and Semantic Web)	géographique	2018	Amalgame de plusieurs licences	oui	rdfs	Pas mentionné
Getty Thesaurus of Geographic Names	Plusieurs équipes (musée, parc nationaux)	historique	2018	Open Data Commons Attribution License v1.0	oui	xml/rdf	Beaucoup de projet ayant pour thème l’histoire
GeoEthno	CNRS, Université Paris Nanterre	Science social	2020	Libre de partager Préciser le que l’on utilise GeoEthno, pas d’utilisation commerciale, pas de modification	Isidore (mais ce n’est pas une API)	rdf skos	Isidore

Dans ce tableau sont présentés uniquement 4 référentiels géographiques. Beaucoup d'autres ont été étudiés mais écartés majoritairement en raison de leur couverture géographique insuffisante.

Le premier référentiel présenté, GeoNames est une base de données géographiques consultable sur internet qui contient plus de 25 millions de noms géographiques et se compose de plus de 11 millions de noms uniques. Ce référentiel est de type wiki avec une équipe qui vérifie les données (pour le mode premium). Plusieurs informations concernant les lieux géographiques sont présentes, comme les coordonnées géographiques ainsi que les différents noms que peuvent porter certaines régions, les données sont mises à jour régulièrement. En ce qui concerne la licence d'utilisation, elle est libre d'accès, il faut uniquement préciser que l'on utilise GeoNames. En ce qui concerne l'accessibilité du référentiel via une API Rest cela est possible, cependant il n'y a pas de Sparql endpoint. Nous pouvons néanmoins télécharger un dump⁶ de GeoNames au format rdf. Ce référentiel est utilisé par de nombreuses entreprises comme Microsoft ou Popfly.

Le second référentiel étudié est LinkedGeoData OpenStreetMap. Ce référentiel est un amalgame de plusieurs projets dont GeoNames, il est développé par AKWS (Agile Knowledge Engineering and semantic Web). En ce qui concerne la licence d'utilisation il faut vérifier chaque projet dont est composé ce référentiel. Le projet n'étant plus mis à jour depuis 2018, certaines données peuvent ne plus être exactes. Une API et un Sparql endpoint existent, mais au moment de la rédaction de ce document les services en ligne sont inutilisables. Le format de données utilisé est le rdfs et aucun projet utilisant ce référentiel n'est mentionné.

Getty Thesaurus of Geographic Names est un référentiel différent des deux précédents, car l'objectif de ce référentiel est plutôt axé sur l'histoire. Cela a un certain intérêt, car de nombreuses villes au cours de l'histoire ont changé de nom ce qui peut être un problème lors de l'analyse de documents. Ce référentiel a été développé par plusieurs musées et parc nationaux. Les données n'ont pas été mises à jour depuis 2018. En ce qui concerne la licence les données sont en accès libre et nous pouvons réutiliser ce projet. Une API est disponible ainsi qu'un Sparql endpoint. Le format des données, correspond à nos besoins et de nombreux projets ont utilisé ce référentiel.

Le dernier référentiel étudié est GeoEthno. Il s'agit d'un projet développé par le CNRS et l'université Paris Nanterre. La thématique de ce référentiel n'est pas orientée géographie, mais science sociale ce qui ne correspond pas vraiment à nos besoins. L'utilisation de ce référentiel est différente des trois autres car il faut utiliser un autre projet (Isidore), de plus lorsqu'on utilise Isidore pour interagir avec Geoethno nous avons accès à GeoNames. Là encore ce référentiel est un amalgame de plusieurs référentiels et les conditions d'utilisation ne correspondent pas à nos besoins.

En conclusion, parmi les critères recherchés, le référentiel géographique répondant le plus à nos besoins est GeoNames. Il possède une couverture mondiale, les ressources sont libres d'accès, de nombreuses données géographiques sont présentes, les données sont disponibles au format du web sémantique. Néanmoins pour pouvoir utiliser les données au format rdf il nous faut télécharger le dump de GeoNames.

⁶ Copie du contenu de GeoNames sur un fichier

3.2.2 Conception et mise en place de la chaîne de traitement

a) Choix d'un outil d'extraction d'entités nommées

Après avoir analysé et sélectionné le référentiel géographique le plus adapté à notre projet, il m'a fallu étudier les différents outils d'extraction d'entités nommées. Là encore il a fallu chercher des critères d'évaluation afin de choisir le meilleur outil.

Les critères sont les suivants :

- Outil Open source
- Multilingue Anglais, Français
- Indication de la Position des entités dans le texte
- Désambiguïsation

Les deux outils testés sont Blink et entity-fishing.

Afin de pouvoir tester ces outils il m'a fallu mettre en place une machine virtuelle pour me permettre d'installer Blink et entity-fishing.

Le premier outil que j'ai testé est Blink, qui est une librairie python de liage d'entités avec Wikipédia. Blink est développé par Facebook research, il utilise une approche basée sur des architectures BERT (c'est un réseau de neurones qui permet de traiter une grande variété de problèmes du traitement automatique de la langue).

Le projet est open source sous licence MIT.

Au travers des différents tests effectué avec cet outil, les résultats sur des textes anglais sont concluants. Toutes les entités nommées sont reconnues et sont liées à un id Wikipédia.

```
des petits ruminants (PR) is a highly contagious disease of small ruminants. The causal agent, PPR Virus (PPRV), is classified into four genetically distinct lineages. Lineage I, originally from Africa, has shown a unique capacity to spread across Africa, the Middle East and Europe. Recent studies have reported its presence in two African countries: Nigeria and Mali. Animals are frequently exchanged between Mali and Nigeria, which could allow the virus to enter and progress in Mali and to other African countries. Here, PPRV samples were collected from sick goats between 2014 and 2017 in both Mali and in Nigeria, on the border with Mali. Partial PPRV nucleoprotein gene was sequenced to identify the genetic lineage of the strains. Our results showed that lineage IV was present in south-eastern Mali in 2017. This is currently the furthest West the lineage has been detected in Africa. Surprisingly, we identified the persistence at least until 2014 of the supposedly extinct lineage I in two regions of Mali, Niger and Nigeria. West PPRV sequences obtained in this study belonged to lineage II, which is dominant in West Africa. Phylogenetic analyses showed a close relationship between sequences obtained at the border between Nigeria and Mali, supporting the hypothesis of an important movement of the virus between the two countries. Understanding the movement of animals between these countries, where the livestock trade is not fully controlled, is very important in the design of efficient control strategies to combat this devastating disease.
```

```
id:110406
title:Pest (organism)
text: A pest is any animal or plant detrimental to humans or human concerns, including crops, livestock and forestry. The term is also used of organisms that cause a nuisance, such as in the home. An older usage is of a deadly epidemic disease, specifically plague.
```

```
id:1085600
title:Ovine rinderpest
text: Ovine rinderpest, also commonly known as peste des petits ruminants (PPR), is a contagious disease primarily affecting goats and sheep; however, camels and wild small ruminants can also be affected. PPR is currently present in North Africa, Central, West and East Asia.
```

```
id:1033309
title:Holin superfamily IV
text: The Holin superfamily IV is a superfamily of integral membrane transport proteins. It is one of the seven different holin superfamilies in total. The Holin superfamily IV includes the TC families: Superfamily IV includes four TC families, which includes TC1, TC2, TC3 and TC4.
```

```
id:72
title:Asia
text: Asia () is Earth's largest and most populous continent, located primarily in the Eastern and Northern Hemispheres. It shares the continental landmass of Eurasia with the continent of Europe and the continental landmass of Afro-Eurasia with both Europe and Africa.
```

```
id:72
title:Asia
text: Asia () is Earth's largest and most populous continent, located primarily in the Eastern and Northern Hemispheres. It shares the continental landmass of Eurasia with the continent of Europe and the continental landmass of Afro-Eurasia with both Europe and Africa.
```

```
id:9453
title:Middle East
text: The Middle East is a transcontinental region centered on Western Asia, Turkey (both Asian and European), and Egypt (which is mostly in North Africa). Saudi Arabia is geographically the largest Middle Eastern nation while Bahrain is the smallest. The course of the Middle East is highly variable.
```

```
id:1036429
title:Africa
text: Africa is the world's second largest and second most-populous continent, being behind Asia in both categories. At about 30.3 million km (11.7 million square miles) including adjacent islands, it covers 6% of Earth's total surface area and 20% of its land area.
```

```
id:3933
title:West Africa
text: West Africa is the westernmost region of Africa. The United Nations defines Western Africa as the 16 countries of Benin, Burkina Faso, Cape Verde, The Gambia, Ghana, Guinea, Guinea-Bissau, Ivory Coast, Liberia, Mali, Mauritania, the Niger, Nigeria, Senegal, Sierra Leone, and Togo.
```

```
id:10496
title:Nigeria
text: Nigeria (), officially the Federal Republic of Nigeria, is a country in West Africa, bordering Niger in the north, Chad in the northeast, Cameroon in the east, and Benin in the west. Its coast in the south is located on the Gulf of Guinea in the Atlantic Ocean.
```

Figure 3 : Résultat Blink avec un texte issu d'agritrop

Le problème majeur de Blink est que la méthode BERT utilisée pour faire fonctionner cet outil a besoin d'un apprentissage en fonction de la langue pour permettre au réseau de neurones de désambiguïser le mieux possible. Cela a pour conséquence de ne pouvoir traiter qu'un seul type de langue alors que les textes fournis par Agritrop sont majoritairement en anglais ou français.

On s'est donc tourné vers un second outil, entity-fishing, développé par la société science-miner qui est aussi la société qui a conçu Grobid, un des logiciels utilisés pour transformer des pdf en document TEI ou txt.

Ce logiciel est open source sous licence Apache 2.

La première différence entre entity-fishing et Blink est la manière dont est fait le liage avec les entités nommées. Blink utilise Wikipedia alors que entity-fishing utilise à la fois Wikipedia, mais aussi Wikidata. De plus, lorsque celui-ci est dans l'incapacité d'effectuer le liage, un « Type » de classe est renvoyé par Grobid ner (name entity recognition). Grobid ner possède 27 types dont un type Location pour les entités géographiques.

Service to call
disambiguate - test

```

{
  "text": "Différents mécanismes d'intervention sont proposés pour transformer les paysages agricoles de manière à ce qu'ils remplissent de multiples fonctions, compatibles avec les objectifs du développement durable. Dans le cas de l'atténuation et de l'adaptation au changement climatique, des politiques incitatives et des mécanismes de rémunération des acteurs locaux sont promues dans le cadre des initiatives REDD+. Sur le terrain, ces interventions visent à assurer simultanément la fourniture de services écosystémiques et le maintien, voire l'amélioration, des moyens d'existence locaux. Dans cet article, nous explorons le rôle que peut jouer l'évaluation participative dans la mise en œuvre de paiements pour services environnementaux au sein de communautés rurales d'Indonésie, du Laos, du Vietnam et de Chine. Engagées dans un processus d'intégration rapide à l'économie de marché, ces communautés ont transformé leurs systèmes d'élevage traditionnels pour s'orienter vers différentes voies d'intensification agricole au cours des dernières années. L'évaluation positive ou négative de ces évolutions fait nécessairement l'objet de jugements de valeur dans les choix effectués entre différentes options possibles. La formulation de ces options en termes de services écosystémiques rendus permet de comparer les trajectoires paysagères et leur impact sur les conditions de vie locales. Sur ces bases, les communautés peuvent explorer des scénarios de transformation de leurs pratiques agricoles, négocier des compromis entre services écosystémiques et identifier les 'gagnants' et les 'perdants' potentiels. Ces simulations ont montré l'importance du calendrier de mise en place des programmes REDD+ par rapport aux transformations agricoles en cours. Elles pointent le risque de s'éloigner des ambitions initiales d'un impact sur les émissions de carbone pour devenir un instrument supplémentaire de développement durable.",
  "entities": {
    "REDD+": {
      "label": "MOT",
      "start": 100,
      "end": 115,
      "type": "MOT"
    },
    "services écosystémiques": {
      "label": "MOT",
      "start": 200,
      "end": 230,
      "type": "MOT"
    },
    "Chine": {
      "label": "MOT",
      "start": 350,
      "end": 365,
      "type": "MOT"
    },
    "Indonésie": {
      "label": "MOT",
      "start": 450,
      "end": 475,
      "type": "MOT"
    },
    "Laos": {
      "label": "MOT",
      "start": 480,
      "end": 495,
      "type": "MOT"
    },
    "Vietnam": {
      "label": "MOT",
      "start": 500,
      "end": 525,
      "type": "MOT"
    }
  },
  "mentions": [
    "REDD+",
    "services écosystémiques",
    "Chine",
    "Indonésie",
    "Laos",
    "Vietnam"
  ]
}

```

WW1
News_1
PubMed_1
PubMed_2
HAL_1
Italiano

News_2
query_1
query_2
query_3
query_4


COVID-19
French
German
Spanish

Annotations
Response

Différents mécanismes d'intervention sont proposés pour transformer les paysages agricoles de manière à ce qu'ils remplissent de multiples fonctions, compatibles avec les objectifs du développement durable. Dans le cas de l'atténuation et de l'adaptation au changement climatique, des politiques incitatives et des mécanismes de rémunération des acteurs locaux sont promues dans le cadre des initiatives REDD+. Sur le terrain, ces interventions visent à assurer simultanément la fourniture de services écosystémiques et le maintien, voire l'amélioration, des moyens d'existence locaux. Dans cet article, nous explorons le rôle que peut jouer l'évaluation participative dans la mise en œuvre de paiements pour services environnementaux au sein de communautés rurales d'Indonésie, du Laos, du Vietnam et de Chine. Engagées dans un processus d'intégration rapide à l'économie de marché, ces communautés ont transformé leurs systèmes d'élevage traditionnels pour s'orienter vers différentes voies d'intensification agricole au cours des dernières années. L'évaluation positive ou négative de ces évolutions fait nécessairement l'objet de jugements de valeur dans les choix effectués entre différentes options possibles. La formulation de ces options en termes de services écosystémiques rendus permet de comparer les trajectoires paysagères et leur impact sur les conditions de vie locales. Sur ces bases, les communautés peuvent explorer des scénarios de transformation de leurs pratiques agricoles, négocier des compromis entre services écosystémiques et identifier les 'gagnants' et les 'perdants' potentiels. Ces simulations ont montré l'importance du calendrier de mise en place des programmes REDD+ par rapport aux transformations agricoles en cours. Elles pointent le risque de s'éloigner des ambitions initiales d'un impact sur les émissions de carbone pour devenir un instrument supplémentaire de développement durable.

CHINE

Normalized: China
Domains: Administration
conf: 0.8222



La Chine, en forme longue la république populaire de Chine (ou République populaire de Chine, RPC, prononcé), parfois appelée Chine populaire, est un pays d'Asie de l'Est. Avec plus d'habitants, soit environ un sixième de la population mondiale, elle est le pays le plus peuplé du monde. Elle compte huit agglomérations de plus de dix millions d'habitants, dont la capitale Pékin, Shanghai, Canton, Shenzhen et Chongqing, ainsi que plus de trente villes d'au moins deux millions d'habitants. Avec selon l'ONU (hors Hong Kong, Macao, et Taïwan) ou de selon The World Factbook, la Chine est également le plus grand pays d'Asie orientale et le troisième ou quatrième plus grand pays du monde par la superficie. La Chine s'étend des côtes de l'océan Pacifique au Pamir et aux Tian Shan, et du désert de Gobi à l'Himalaya et au nord de la péninsule indochinoise.

Wikidata statements

References: W I

Figure 4: Résultat d'entity-fishing avec un texte issu d'Agritrop

L'outil fonctionne avec plusieurs langues (Français, anglais, italien, allemand, espagnol), il n'utilise pas la méthode BERT pour fonctionner, ce qui est un avantage, car il n'y aura pas besoin de réentraîner un modèle à chaque changement de langue.

Entity-fishing possède un service en ligne qui permet de faire des tests avant de l'installer sur la machine virtuelle du projet.

Cela m'a permis d'analyser les premiers résultats et de les comparer à ceux obtenus avec Blink.

Entity-fishing renvoie l'ensemble des entités nommées sous forme d'un json avec les positions dans le texte de chaque concept.

Cet outil répondant à tous nos critères de sélection nous avons décidé de l'utiliser pour la suite du développement.

b) Extraction des entités nommées géographiques et alignement avec GeoNames

En sortie d'entity-fishing, on récupère une liste d'entités nommées de tout type : entités nommées généralistes et géographiques. L'un des objectifs est d'isoler les entités géographiques des autres.

Pour pouvoir réaliser cette chaîne de traitement j'ai travaillé en plusieurs étapes :

- Travail technique : récupérer et traiter le fichier json en sortie d'entity-fishing.
- Isoler et traiter les entités nommées géographiques pour les aligner avec GeoNames.

Travail technique

Pour récupérer et traiter le fichier json obtenu en sortie d'entity-fishing, il existe une librairie utilisable en python. Cette librairie développée par Hirneios avec l'aide de l'équipe de science-miner et sous licence Apache 2.0 m'a permis de communiquer avec entity-fishing installé sur la machine virtuelle.

J'ai utilisé putty afin d'établir la connexion entre mon ordinateur et la machine virtuelle où est installé entity-fishing. Pour fonctionner, la librairie avait seulement besoin d'un point d'accès vers le port de la machine virtuelle associé à entity-fishing.

À la suite de cet ajustement technique, nous avons décidé de construire la chaîne de traitement en python afin de pouvoir utiliser cette librairie mais aussi car python est très utilisé pour effectuer du traitement automatique du langage

Traitement des entités nommées géographiques

Le résultat en sortie de entity-fishing est au format json et est composé de trois parties :

- Des informations concernant la version de l'outil : le texte envoyé, la date et heure du run ainsi que la langue associée au texte
- Une partie avec des « global_categories » associées à Wikipedia
- Une partie contenant toutes les entités nommées

```
{
  "software": "entity-fishing",
  "version": "0.0.4",
  "date": "2021-08-21T13:47:28.743Z",
  "runtime": 1754,
  "nbest": false,
  "text": "Edmond Ludlow (vers 1617-1692) est un parlementaire anglais, plus connu pour son implication dans l'exécution de Charles Ier, et pour ses mémoires, publiés à titre posthume et qui sont devenus une source importante pour les historiens des Guerres des Trois Royaumes. Après avoir servi dans les guerres civiles anglaises, Ludlow a été élu membre du Long Parlement. Après la création du Commonwealth en 1649, il est nommé adjoint de Ireton, commandant des forces du Parlement en Irlande, avant de rompre avec Oliver Cromwell lors de la création du Protectorat. Après la Restauration, Ludlow part en exil en Suisse, où il passe une grande partie du reste de sa vie.",
  "language": {
    "lang": "fr",
    "conf": 0
  },
  "global_categories": [
    {
      "weight": 0.007751937984496138,
      "source": "wikipedia-fr",
      "category": "Langue de Saint-Christophe-et-Nièves",
      "page_id": 553481
    },
    {
      "weight": 0.007751937984496138,
      "source": "wikipedia-fr",
      "category": "Langue de Papouasie-Nouvelle-Guinée",
      "page_id": 547084
    },
    {
      "weight": 0.007751937984496138,
      "source": "wikipedia-fr",
      "category": "Langue de Sainte-Lucie",
      "page_id": 553484
    }
  ],
  "entities": [
    {
      "rawName": "Edmond Ludlow",
      "offsetStart": 0,
      "offsetEnd": 13,
      "confidence_score": 0.9921,
      "wikipediaExternalRef": 8331104,
      "wikidataId": "Q5339650",
      "domains": [
        "Administration"
      ]
    }
  ]
}
```

Figure 5 : Exemple de réponse au format json de entity-fishing

Plus précisément, voici les informations d'importance pour nous :

- Au début du fichier, seule la langue nous intéresse pour pouvoir différencier correctement les textes anglais et français.
- Les « global categories » ne nous intéressent pas, car dans cette tâche on se concentre uniquement sur les entités nommées géographiques et celles-ci ne comportent aucune information utile de ce genre.
- La partie « entities » est la plus importante : elle contient toutes les entités nommées (géographiques et autres) ainsi que divers renseignements.

On peut séparer ces entités en deux parties, celles qui possèdent un id Wikidata et Wikipédia et celles qui possèdent seulement un type sans id.

Les entités nommées qui possèdent les deux id sont celles qui ont pu être liées avec Wikipedia.

Les entités qui possèdent un type NER sont celles que entity-fishing n'a pas réussi à désambigüiser. Ce type est obtenu grâce à grobid tel qu'expliqué précédemment. Dans notre cas on s'intéressera uniquement à celles qui ont le type « LOCATION »

```
{
  "rawName": "Irlande",
  "type": "LOCATION",
  "offsetStart": 477,
  "offsetEnd": 484,
  "confidence_score": 0
},
```

Figure 6 : Exemple d'une entité nommée géographique non désambigüisée

Chaque entité correctement désambigüisée est composée de :

- Son « rawName »
- De sa position dans le texte (début et fin)
- D'un score de confiance par rapport à la désambigüisation, confidence_score
- De deux id (Wikidata, Wikipedia)
- D'un ou plusieurs « domains »

```
{
  "rawName": "services écosystémiques",
  "offsetStart": 493,
  "offsetEnd": 516,
  "confidence_score": 0.9948,
  "wikipediaExternalRef": 858238,
  "wikidataId": "Q295865",
  "domains": [
    "Environment",
    "Hydraulics",
    "Plants"
  ]
},
```

Figure 7 : Exemple d'une entité nommée désambigüisée

En sortie de traitement j'obtiens donc une liste d'entités nommées :

- Soit désambigüisée avec des id Wikipedia, Wikidata
- Soit non désambigüisée avec un type NER

Travail préalable technique : reformater les résultats pour pouvoir les exploiter

- Création d'un programme python qui prend en entrée un répertoire contenant des fichiers texte avec les résumés de quatre documents issus d'Agritrop et qui renvoie en sortie un autre répertoire qui va contenir les quatre fichiers json.

Ce programme python contient plusieurs fonctions :

- Une fonction read qui prend en entrée un répertoire et un nom de fichier qui va me permettre de lire chaque fichier txt (en utf-8)
- Une fonction pour éliminer les parties inutiles et garder seulement la métadonnée sur la langue du texte et toutes les entités
- Une fonction d'écriture pour écrire le résultat dans un fichier json dans le répertoire de sortie en gardant un encoding utf-8 et en gardant le même nom de fichier que celui des txt, car chaque fichier txt se nomme avec id Agritrop afin de pouvoir repérer plus facilement le document concerné.
- Deux fonctions en lien avec la librairie python d'entity-fishing afin de communiquer avec l'outil
- Un programme principal « main » dans lequel on va parcourir chaque fichier du répertoire d'entrée pour obtenir un résultat dans le répertoire de sortie.

À ce stade, on a donc récupéré toutes les entités nommées et reformaté le résultat.

Problématique : Comment isoler les entités nommées géographiques des autres et les aligner avec GeoNames ?

Une première idée était d'utiliser les « domains » associés à chaque entité afin de répertorier celles qui sont géographiques, mais cela était impossible, car les « domains » n'étaient pas pertinents. Du coup, cette partie a été retirée du résultat.

Les seuls éléments pouvant nous aider restaient les id Wikidata et Wikipedia. En ce qui concerne les entités avec un type location, elles seront traitées ultérieurement.

Pour chaque EN géographique trouvée, l'id Wikipedia renvoie sur une page qui ne contient pas de lien vers GeoNames.

En revanche, on peut trouver un pointeur vers GeoNames dans Wikidata, ce qui me permet de lier les EN extraites avec GeoNames.

Travail technique pour séparer les EN géographiques des autres et les lier avec GeoNames.

Pour cela j'ai créé un second programme python afin d'effectuer des requêtes sparql sur le sparql endpoint de Wikidata de manière à récupérer les id GeoNames grâce aux id Wikidata.

La requête sparql pour interroger Wikidata n'est pas intégrée au code python, mais séparée dans un fichier .txt (bonne pratique).

Afin de pouvoir récupérer le fichier txt contenant ma requête sparql, j'utilise la librairie string puis pour interroger Wikidata, la librairie SPARQLWrapper.

Ce programme python contient les fonctions suivantes :

- Une fonction read_json pour lire le json obtenu à l'étape précédente
- Une fonction read_txt pour lire la requête Sparql
- Une fonction write afin d'écrire le résultat dans un nouveau json
- Une fonction Get_value_id_GeoNames_wrapper afin d'effectuer la requête à Wikidata
- Une fonction Traitement_Wikidata afin de récupérer le résultat de la requête sparql et de récupérer uniquement l'id GeoNames
- Un programme principal « main » pour exécuter l'ensemble.

La chaîne de traitement est volontairement divisée en deux programmes afin de pouvoir analyser chaque étape et de pouvoir effectuer des analyses sur les résultats.

Lorsqu'un id Wikidata possède un id GeoNames, on l'associe avec l'entité nommée auquel il appartient. Il est ensuite gardé et écrit dans un nouveau json situé dans un nouveau répertoire. Sinon il est supprimé de la liste des entités. En ce qui concerne les entités géographiques avec un type LOCATION, elles sont gardées dans le json sans modification.

Une fois le programme terminé, nous obtenons un json avec uniquement des entités géographiques, soit associées à un id GeoNames, soit à un type LOCATION.

Après avoir optimisé certaines parties, j'ai continué à explorer diverses solutions pour trouver un moyen d'arriver à lier les entités qui possèdent seulement un type LOCATION à une base de connaissance.

Pour pouvoir résoudre ce problème, on renvoie sans contexte de phrase les entités géographiques dans entity-fishing pour les désambiguïser avec Wikidata.

Cependant il y a deux inconvénients :

- Le premier est que le « confidence score » a moins de valeur car l'EN est désambiguïlée sans contexte
- Le second est que le risque de faux positifs est augmenté, pour cela un score de confiance d'au moins 10% est requis sinon l'entité géographique est éliminée. même si l'outil la déclare comme entité géographique

c) Analyse des résultats

Les résultats que je vais maintenant exposer ont été obtenus après traitement de deux articles scientifiques issus de l'archive ouverte Agritrop.

Le premier est un article intitulé « Persistence of the historical lineage I of West Africa against the ongoing spread of the Asian lineage of peste des petits ruminants virus » écrit en langue anglaise (id dans Agritrop : 598198).

Cet article possède 36 entités géographiques distinctes identifiées manuellement par un documentaliste et chacune d'entre elles peut être répétée plusieurs fois dans le texte (124 fois avec la répétition).

Sur ces 36 entités distinctes, le programme en trouve 26/36 avec deux entités qui sont des faux positifs, soit un pourcentage de réussite de 63%.

Ce pourcentage de réussite est à nuancer, car certaines des entités géographiques présentes dans le texte ne possèdent pas de page Wikipedia ni Wikidata (il s'agit de petits villages de l'ouest de l'Afrique). Par conséquent il est impossible d'arriver à désambiguïser ou à lier ces entités avec entity-fishing. Si l'on s'en tient aux entités nommées (EN) qui ont un id Wikipédia, alors le taux de réussite est meilleur.

Par exemple, dans cet article, sept EN ne peuvent être trouvées avec la chaîne de traitement car elles ne sont pas répertoriées dans Wikipédia. En excluant ces sept EN, on obtient alors un pourcentage de réussite de 77%.

Les 3 EN typées LOCATION n'ont pas été **retenu** car non répertoriées dans GeoNames.

Dans le tableau ci-dessous, on peut voir les EN de la publication avec dans la première colonne le label d'EN, dans la seconde le nombre de répétitions et dans la troisième colonne une note explicative.

Africa	2 sur 2	
Asia	1 sur 1	
Bamako	1 sur 1	
Benin	2 sur 2	
Burkina Faso	3 sur 3	
Central, East, North and West Africa	0 sur 1	impossible car entité géographique composé
Dialan	0 sur 2	n'existe pas dans la base de donnée de wikidata ni geonames
East	1 fois	Est de l'afrique mais problème de désambiguisation
East Africa	2 sur 2	
France	0 sur 1	Problème même lors de la désambiguisation dans le pdf
Ghana	1 sur 1	
Ivory Coast	1 sur 1	
Kayes	3 sur 2	
Kayes region	3 sur 4	
Kedougou	3 sur 3	
Kenieba	1 sur 1	
Kolondieba	1 sur 1	
Kopropin	0 sur 1	n'existe pas dans la base de donnée de wikidata ni geonames
Krounikoto	0 sur 1	n'existe pas dans la base de donnée de wikidata ni geonames
Mali	29 sur 29	
Montpellier	0 sur 1	Problème même lors de la désambiguisation dans le pdf
Mauritania	1 sur 1	
Middle East	1 sur 1	
Mopti	3 sur 4	
Niger	11 sur 11	
Nigeria	7 sur 7	
Samako	0 sur 1	n'existe pas dans la base de donnée de wikidata ni geonames
Segou	5 sur 5	
Senegal	14 sur 14	
Seroume	0 sur 1	n'existe pas dans la base de donnée de wikidata ni geonames
Sikasso	1 sur 1	
Sitakili	0 sur 1	Problème dans le fichier txt car reconnu dans le pdf
Tambacounda	1 sur 1	
united kingdom	1 sur 1	
west	1 fois	ouest de l'afrique mais problème de désambiguisation
West and East Africa	0 sur 1	impossible car entité géographique composé
West Africa	8 sur 8	
West african	6 sur 6	associe le même id que west africa
Article 598198		

Figure 8 : Analyse des résultats pour l'article 598198

À partir de l'analyse des résultats présentés précédemment, nous pouvons calculer les scores de précision⁷ et de rappel⁸.

598198	Entités nommées présent	Entités nommées absent
Réponse du test positive	Vrai positif : 111	Faux positif : 3
Réponse du test négative	Faux négatif : 10	Vrai négatif : /

- Précision : $111 / (111 + 3) = 97\%$
- Rappel : $111 / (111 + 10) = 91\%$

⁷ La précision donne le pourcentage de réponses correctes.

⁸ Le rappel donne le pourcentage des réponses correctes qui sont données.

Le deuxième est un article intitulé « Explorer l'impact environnemental des transformations agraires en Asie du Sud-Est grâce à l'évaluation participative des services écosystémiques » écrit en français (id dans Agritrop : 597393).

Cet article possède 22 entités géographiques distinctes identifiées manuellement par un documentaliste et chacune d'entre elles peut être répétée plusieurs fois dans le texte (52 fois avec la répétition).

Sur ces 22 entités distinctes, le programme en trouve 13/22 soit un pourcentage de réussite de 59%.

Ce pourcentage de réussite et aussi à nuancer car certaines entités géographiques ne sont pas présentes dans nos différents référentiels.

3 entités géographiques ne peuvent être trouvées, alors le taux de réussite est de 68%.

6 entités étaient typées LOCATION n'ont pas été retenues car non répertoriées dans GeoNames.

Dans le tableau ci-dessous, on peut voir les EN de la publication avec dans la première colonne le label d'EN, dans la seconde le nombre de répétitions et dans la troisième colonne une note explicative.

Afrique	1 sur 1	
Amérique Latine	1 sur 1	
Asie	1 sur 1	
Asie du Sud-Est	2 sur 3	
Chine	7 sur 8	
Chinois	3 sur 3	
district de Con Cuong	0 sur 1	Problème même lors de la désambiguïsation dans le pdf
district de Hiem	0 sur 1	n'existe pas dans la base de données de Wikidata ni Geonames
district de Kutai Barat	0 sur 1	Problème même lors de la désambiguïsation dans le pdf
Huaphan	1 sur 1	
Indonésie	3 sur 3	
Indonésiens	0 sur 3	Problème même lors de la désambiguïsation dans le pdf
Kalimantan	0 sur 1	Problème dans le fichier txt car reconnu dans le pdf
Kalimantan Est	0 sur 1	Problème dans le fichier mais ne trouve pas Kalimantan Est, l'outil fait le lien avec Kalimantan
Laos	7 sur 7	
Laotiens	0 sur 2	Problème même lors de la désambiguïsation dans le pdf
préfecture de Xishuangbanna	0 sur 1	Problème dans le fichier txt car reconnu dans le pdf
province de Nghe An	1 sur 1	
province du Yunnan	1 sur 1	
région chinoise du Xishuangbanna	0 sur 1	Problème dans le fichier txt car reconnu dans le pdf
Thaïlande	1 sur 1	
Vietnam	9 sur 9	
Article 597393		

Figure 9 : Analyse des résultats pour l'article 597393

À partir de l'analyse des résultats présentés précédemment, nous pouvons à nouveau calculer les scores de précision et de rappel pour ce second article.

- Précision : $38 / (38 + 0) = 100\%$
- Rappel : $38 / (38 + 14) = 73\%$

5597393	Entités nommées présent	Entités nommées absent
Réponse du test positive	Vrai positif : 38	Faux positif : 0
Réponse du test négative	Faux négatif : 14	Vrai négatif : /

Le score de rappel obtenu précédemment pour les documents ne prennent pas en compte le fait que certaines entités nommées géographiques ne possèdent pas d'id dans Wikipedia, Wikidata et GeoNames.

Si l'on refait les calculs pour les 2 documents en enlevant en enlevant ces entités nommées là on obtient :

Pour le premier document, nous avons 119 entités nommées géographiques en enlevant celles citées précédemment.

598198	Entités nommées présente	Entités nommées absente
Réponse du test positive	Vrai positif : 111	Faux positif : 3
Réponse du test négative	Faux négatif : 5	Vrai négatif : /

- Rappel : $111 / (111 + 5) = 95\%$

Et pour le second document, nous avons 46 entités nommées géographiques en enlevant celle cité précédemment.

599793	Entités nommées présente	Entités nommées absente
Réponse du test positive	Vrai positif : 38	Faux positif : 0
Réponse du test négative	Faux négatif : 8	Vrai négatif : /

- Rappel : $38 / (38 + 8) = 82\%$

En conclusion, sur les 2 documents analysés nous obtenons un score de précision de plus de 95% nous permettant de voir que le nombre de faux positifs est extrêmement bas.

En ce qui concerne le score de rappel, nous avons aussi de très bons résultats, dans les 2 cas supérieurs à 70 %.

Lorsqu'on enlève les entités nommées qui ne se trouvent dans aucune des bases de connaissance (Wikipédia, Wikidata, GeoNames), le score de rappel est supérieur à 80% soit une augmentation de 10%.

La chaîne de traitement mise en place obtient un score de plus de 75% de réussite dans le meilleur des cas.

Le texte en anglais possède un meilleur score que celui en français, probablement parce que les outils utilisés fonctionnent mieux en langue anglaise.

3.2.3 Enrichissement des données et visualisation cartographique

Dans cette section, on explique les procédés mis en place pour visualiser les EN géographiques sur une carte.

Pour pouvoir enrichir les données obtenues en sortie de la chaîne de traitement, on utilise le dump de GeoNames. Le premier problème rencontré est que ce dump n'est pas un fichier rdf valide (voir figure ci-dessous). Il s'agit d'une concaténation de plusieurs milliards de fichier rdf.

Ce fichier contient des éléments associés à l'id GeoNames comme les coordonnées géographiques, les différents noms d'un même endroit, s'il s'agit d'un **contient**, d'un pays ou d'une ville.

```
1https://sws.geonames.org/3/
2<?xml version="1.0" encoding="UTF-8" standalone="no"?><rdf:RDF xmlns:cc="http://creativecommons.org/ns#" xmlns:dcterms="http://purl.org/dc/terms/" xmlns:foaf="http://xmlns.com/foaf/0.1/" xmlns:gn="http://www.geonames.org/ont
3https://sws.geonames.org/4/
4<?xml version="1.0" encoding="UTF-8" standalone="no"?><rdf:RDF xmlns:cc="http://creativecommons.org/ns#" xmlns:dcterms="http://purl.org/dc/terms/" xmlns:foaf="http://xmlns.com/foaf/0.1/" xmlns:gn="http://www.geonames.org/ont
5https://sws.geonames.org/5/
6<?xml version="1.0" encoding="UTF-8" standalone="no"?><rdf:RDF xmlns:cc="http://creativecommons.org/ns#" xmlns:dcterms="http://purl.org/dc/terms/" xmlns:foaf="http://xmlns.com/foaf/0.1/" xmlns:gn="http://www.geonames.org/ont
7https://sws.geonames.org/6/
8<?xml version="1.0" encoding="UTF-8" standalone="no"?><rdf:RDF xmlns:cc="http://creativecommons.org/ns#" xmlns:dcterms="http://purl.org/dc/terms/" xmlns:foaf="http://xmlns.com/foaf/0.1/" xmlns:gn="http://www.geonames.org/ont
9https://sws.geonames.org/7/
10<?xml version="1.0" encoding="UTF-8" standalone="no"?><rdf:RDF xmlns:cc="http://creativecommons.org/ns#" xmlns:dcterms="http://purl.org/dc/terms/" xmlns:foaf="http://xmlns.com/foaf/0.1/" xmlns:gn="http://www.geonames.org/ont
11https://sws.geonames.org/8/
12<?xml version="1.0" encoding="UTF-8" standalone="no"?><rdf:RDF xmlns:cc="http://creativecommons.org/ns#" xmlns:dcterms="http://purl.org/dc/terms/" xmlns:foaf="http://xmlns.com/foaf/0.1/" xmlns:gn="http://www.geonames.org/ont
13https://sws.geonames.org/9/
14<?xml version="1.0" encoding="UTF-8" standalone="no"?><rdf:RDF xmlns:cc="http://creativecommons.org/ns#" xmlns:dcterms="http://purl.org/dc/terms/" xmlns:foaf="http://xmlns.com/foaf/0.1/" xmlns:gn="http://www.geonames.org/ont
15https://sws.geonames.org/10/
16<?xml version="1.0" encoding="UTF-8" standalone="no"?><rdf:RDF xmlns:cc="http://creativecommons.org/ns#" xmlns:dcterms="http://purl.org/dc/terms/" xmlns:foaf="http://xmlns.com/foaf/0.1/" xmlns:gn="http://www.geonames.org/ont
17https://sws.geonames.org/11/
18<?xml version="1.0" encoding="UTF-8" standalone="no"?><rdf:RDF xmlns:cc="http://creativecommons.org/ns#" xmlns:dcterms="http://purl.org/dc/terms/" xmlns:foaf="http://xmlns.com/foaf/0.1/" xmlns:gn="http://www.geonames.org/ont
19https://sws.geonames.org/12/
20<?xml version="1.0" encoding="UTF-8" standalone="no"?><rdf:RDF xmlns:cc="http://creativecommons.org/ns#" xmlns:dcterms="http://purl.org/dc/terms/" xmlns:foaf="http://xmlns.com/foaf/0.1/" xmlns:gn="http://www.geonames.org/ont
21https://sws.geonames.org/13/
22<?xml version="1.0" encoding="UTF-8" standalone="no"?><rdf:RDF xmlns:cc="http://creativecommons.org/ns#" xmlns:dcterms="http://purl.org/dc/terms/" xmlns:foaf="http://xmlns.com/foaf/0.1/" xmlns:gn="http://www.geonames.org/ont
23https://sws.geonames.org/14/
24<?xml version="1.0" encoding="UTF-8" standalone="no"?><rdf:RDF xmlns:cc="http://creativecommons.org/ns#" xmlns:dcterms="http://purl.org/dc/terms/" xmlns:foaf="http://xmlns.com/foaf/0.1/" xmlns:gn="http://www.geonames.org/ont
25https://sws.geonames.org/15/
26<?xml version="1.0" encoding="UTF-8" standalone="no"?><rdf:RDF xmlns:cc="http://creativecommons.org/ns#" xmlns:dcterms="http://purl.org/dc/terms/" xmlns:foaf="http://xmlns.com/foaf/0.1/" xmlns:gn="http://www.geonames.org/ont
27https://sws.geonames.org/16/
28<?xml version="1.0" encoding="UTF-8" standalone="no"?><rdf:RDF xmlns:cc="http://creativecommons.org/ns#" xmlns:dcterms="http://purl.org/dc/terms/" xmlns:foaf="http://xmlns.com/foaf/0.1/" xmlns:gn="http://www.geonames.org/ont
29https://sws.geonames.org/17/
30<?xml version="1.0" encoding="UTF-8" standalone="no"?><rdf:RDF xmlns:cc="http://creativecommons.org/ns#" xmlns:dcterms="http://purl.org/dc/terms/" xmlns:foaf="http://xmlns.com/foaf/0.1/" xmlns:gn="http://www.geonames.org/ont
31https://sws.geonames.org/18/
32<?xml version="1.0" encoding="UTF-8" standalone="no"?><rdf:RDF xmlns:cc="http://creativecommons.org/ns#" xmlns:dcterms="http://purl.org/dc/terms/" xmlns:foaf="http://xmlns.com/foaf/0.1/" xmlns:gn="http://www.geonames.org/ont
33https://sws.geonames.org/19/
34<?xml version="1.0" encoding="UTF-8" standalone="no"?><rdf:RDF xmlns:cc="http://creativecommons.org/ns#" xmlns:dcterms="http://purl.org/dc/terms/" xmlns:foaf="http://xmlns.com/foaf/0.1/" xmlns:gn="http://www.geonames.org/ont
35https://sws.geonames.org/20/
36<?xml version="1.0" encoding="UTF-8" standalone="no"?><rdf:RDF xmlns:cc="http://creativecommons.org/ns#" xmlns:dcterms="http://purl.org/dc/terms/" xmlns:foaf="http://xmlns.com/foaf/0.1/" xmlns:gn="http://www.geonames.org/ont
37https://sws.geonames.org/21/
38<?xml version="1.0" encoding="UTF-8" standalone="no"?><rdf:RDF xmlns:cc="http://creativecommons.org/ns#" xmlns:dcterms="http://purl.org/dc/terms/" xmlns:foaf="http://xmlns.com/foaf/0.1/" xmlns:gn="http://www.geonames.org/ont
39https://sws.geonames.org/22/
40<?xml version="1.0" encoding="UTF-8" standalone="no"?><rdf:RDF xmlns:cc="http://creativecommons.org/ns#" xmlns:dcterms="http://purl.org/dc/terms/" xmlns:foaf="http://xmlns.com/foaf/0.1/" xmlns:gn="http://www.geonames.org/ont
41https://sws.geonames.org/23/
42<?xml version="1.0" encoding="UTF-8" standalone="no"?><rdf:RDF xmlns:cc="http://creativecommons.org/ns#" xmlns:dcterms="http://purl.org/dc/terms/" xmlns:foaf="http://xmlns.com/foaf/0.1/" xmlns:gn="http://www.geonames.org/ont
43https://sws.geonames.org/24/
44<?xml version="1.0" encoding="UTF-8" standalone="no"?><rdf:RDF xmlns:cc="http://creativecommons.org/ns#" xmlns:dcterms="http://purl.org/dc/terms/" xmlns:foaf="http://xmlns.com/foaf/0.1/" xmlns:gn="http://www.geonames.org/ont
45https://sws.geonames.org/25/
46<?xml version="1.0" encoding="UTF-8" standalone="no"?><rdf:RDF xmlns:cc="http://creativecommons.org/ns#" xmlns:dcterms="http://purl.org/dc/terms/" xmlns:foaf="http://xmlns.com/foaf/0.1/" xmlns:gn="http://www.geonames.org/ont
47https://sws.geonames.org/26/
48<?xml version="1.0" encoding="UTF-8" standalone="no"?><rdf:RDF xmlns:cc="http://creativecommons.org/ns#" xmlns:dcterms="http://purl.org/dc/terms/" xmlns:foaf="http://xmlns.com/foaf/0.1/" xmlns:gn="http://www.geonames.org/ont
49https://sws.geonames.org/27/
50<?xml version="1.0" encoding="UTF-8" standalone="no"?><rdf:RDF xmlns:cc="http://creativecommons.org/ns#" xmlns:dcterms="http://purl.org/dc/terms/" xmlns:foaf="http://xmlns.com/foaf/0.1/" xmlns:gn="http://www.geonames.org/ont
51https://sws.geonames.org/28/
52<?xml version="1.0" encoding="UTF-8" standalone="no"?><rdf:RDF xmlns:cc="http://creativecommons.org/ns#" xmlns:dcterms="http://purl.org/dc/terms/" xmlns:foaf="http://xmlns.com/foaf/0.1/" xmlns:gn="http://www.geonames.org/ont
53https://sws.geonames.org/29/
54<?xml version="1.0" encoding="UTF-8" standalone="no"?><rdf:RDF xmlns:cc="http://creativecommons.org/ns#" xmlns:dcterms="http://purl.org/dc/terms/" xmlns:foaf="http://xmlns.com/foaf/0.1/" xmlns:gn="http://www.geonames.org/ont
55https://sws.geonames.org/30/
56<?xml version="1.0" encoding="UTF-8" standalone="no"?><rdf:RDF xmlns:cc="http://creativecommons.org/ns#" xmlns:dcterms="http://purl.org/dc/terms/" xmlns:foaf="http://xmlns.com/foaf/0.1/" xmlns:gn="http://www.geonames.org/ont
57https://sws.geonames.org/31/
58<?xml version="1.0" encoding="UTF-8" standalone="no"?><rdf:RDF xmlns:cc="http://creativecommons.org/ns#" xmlns:dcterms="http://purl.org/dc/terms/" xmlns:foaf="http://xmlns.com/foaf/0.1/" xmlns:gn="http://www.geonames.org/ont
59https://sws.geonames.org/32/
60<?xml version="1.0" encoding="UTF-8" standalone="no"?><rdf:RDF xmlns:cc="http://creativecommons.org/ns#" xmlns:dcterms="http://purl.org/dc/terms/" xmlns:foaf="http://xmlns.com/foaf/0.1/" xmlns:gn="http://www.geonames.org/ont
61https://sws.geonames.org/33/
62<?xml version="1.0" encoding="UTF-8" standalone="no"?><rdf:RDF xmlns:cc="http://creativecommons.org/ns#" xmlns:dcterms="http://purl.org/dc/terms/" xmlns:foaf="http://xmlns.com/foaf/0.1/" xmlns:gn="http://www.geonames.org/ont
63https://sws.geonames.org/34/
```

Figure 10 : Dump de GeoNames avant mise en conformité

Travail technique

Pour résoudre ce problème, j'ai créé un troisième fichier python permettant de modifier le fichier rdf et de le rendre conforme. Pour cela, j'ai effectué des traitements en utilisant les expressions régulières.

Dans la figure ci-dessous on montre le résultat obtenu.

```

1 <?xml version="1.0" encoding="UTF-8" standalone="no"?>
2 <rdf:RDF xmlns:cc="http://creativecommons.org/ns#" xmlns:dcterms="http://purl.org/dc/terms/" xmlns:foaf="http://xmlns.com/foaf/0.1/" xmlns:gn="http://www.geonames.org/ontology#" xmlns:owl="http://www.w3.org/2002/07/owl#"
3 <gn:Feature rdf:about="https://sws.geonames.org/3/">
4   <rdfs:isDefinedBy rdf:resource="https://sws.geonames.org/3/about.rdf"/>
5   <gn:name>Zamin Sukhteh</gn:name>
6   <gn:featureClass rdf:resource="https://www.geonames.org/ontology#S"/>
7   <gn:featureCode rdf:resource="https://www.geonames.org/ontology#S.CRRL"/>
8   <gn:countryCode>IR</gn:countryCode>
9   <wgs84_pos:lat>32.45831</wgs84_pos:lat>
10  <wgs84_pos:long>48.96335</wgs84_pos:long>
11  <gn:parentFeature rdf:resource="https://sws.geonames.org/3262991"/>
12  <gn:parentCountry rdf:resource="https://sws.geonames.org/130758"/>
13  <gn:parentADM1 rdf:resource="https://sws.geonames.org/127082"/>
14  <gn:nearbyFeatures rdf:resource="https://sws.geonames.org/3/nearby.rdf"/>
15  <gn:locationMap rdf:resource="https://www.geonames.org/3/zamin-sukhteh.html"/>
16 </gn:Feature>
17
18 <gn:Feature rdf:about="https://sws.geonames.org/4/">
19   <rdfs:isDefinedBy rdf:resource="https://sws.geonames.org/4/about.rdf"/>
20   <gn:name>Rudkhaneh-ye Ab-e Zalek</gn:name>
21   <gn:alternateName>Rudkhaneh-ye Ab-e Zalek</gn:alternateName>
22   <gn:alternateName>Rudkhaneh-ye Ab-e Zaleki</gn:alternateName>
23   <gn:featureClass rdf:resource="https://www.geonames.org/ontology#H"/>
24   <gn:featureCode rdf:resource="https://www.geonames.org/ontology#H.STM"/>
25   <gn:countryCode>IR</gn:countryCode>
26   <wgs84_pos:lat>32.93273</wgs84_pos:lat>
27   <wgs84_pos:long>48.76505</wgs84_pos:long>
28   <gn:parentFeature rdf:resource="https://sws.geonames.org/127082"/>
29   <gn:parentCountry rdf:resource="https://sws.geonames.org/130758"/>
30   <gn:parentADM1 rdf:resource="https://sws.geonames.org/127082"/>
31   <gn:nearbyFeatures rdf:resource="https://sws.geonames.org/4/nearby.rdf"/>
32   <gn:locationMap rdf:resource="https://www.geonames.org/4/rudkhaneh-ye-ab-e-zalek.html"/>
33 </gn:Feature>
34
35 <gn:Feature rdf:about="https://sws.geonames.org/5/">
36   <rdfs:isDefinedBy rdf:resource="https://sws.geonames.org/5/about.rdf"/>
37   <gn:name>Yekahi</gn:name>
38   <gn:featureClass rdf:resource="https://www.geonames.org/ontology#P"/>
39   <gn:featureCode rdf:resource="https://www.geonames.org/ontology#P.PPL"/>
40   <gn:countryCode>IR</gn:countryCode>
41   <wgs84_pos:lat>32.5</wgs84_pos:lat>
42   <wgs84_pos:long>48.9</wgs84_pos:long>
43   <gn:parentFeature rdf:resource="https://sws.geonames.org/127082"/>
44   <gn:parentCountry rdf:resource="https://sws.geonames.org/130758"/>
45   <gn:parentADM1 rdf:resource="https://sws.geonames.org/127082"/>
46   <gn:nearbyFeatures rdf:resource="https://sws.geonames.org/5/nearby.rdf"/>
47   <gn:locationMap rdf:resource="https://www.geonames.org/5/yekahi.html"/>
48 </gn:Feature>
49
50 <gn:Feature rdf:about="https://sws.geonames.org/6/">
51   <rdfs:isDefinedBy rdf:resource="https://sws.geonames.org/6/about.rdf"/>
52   <gn:name>Ab-e Yas</gn:name>
53   <gn:featureClass rdf:resource="https://www.geonames.org/ontology#H"/>
54   <gn:featureCode rdf:resource="https://www.geonames.org/ontology#H.STM"/>
55   <gn:countryCode>IR</gn:countryCode>
56   <wgs84_pos:lat>32.8</wgs84_pos:lat>
57   <wgs84_pos:long>48.8</wgs84_pos:long>
58   <gn:parentFeature rdf:resource="https://sws.geonames.org/127082"/>
59   <gn:parentCountry rdf:resource="https://sws.geonames.org/130758"/>
60   <gn:parentADM1 rdf:resource="https://sws.geonames.org/127082"/>
61   <gn:nearbyFeatures rdf:resource="https://sws.geonames.org/6/nearby.rdf"/>
62   <gn:locationMap rdf:resource="https://www.geonames.org/6/ab-e-yas.html"/>
63 </gn:Feature>

```

Figure 11 : Dump de GeoNames après mise en conformité

À ce stade, on possède deux sources de données à exploiter : le fichier json obtenu après les traitements en sortie de entity-fishing et le fichier reformaté du dump GeoNames. Dans la section suivante on va expliquer comment lier les deux et exploiter leurs informations.

Travail technique

Un triplestore Virtuoso a été déployé sur la VM du projet de manière à : stocker les données et déployer un Sparql endpoint permettant de les interroger .

Une fois l'installation terminée, j'ai pu commencer à créer une requête sparql afin de pouvoir récupérer l'ensemble des données en fonction d'un id GeoNames.

Ces données (surtout les coordonnées géographiques) vont me permettre de créer une visualisation cartographique pour situer correctement toutes les entités géographiques d'un document.

Pour réussir cette visualisation, je me suis appuyé sur mes connaissances acquises lors de mon premier semestre : librairie React javascript, requêtes Sparql.

Travail technique

J'ai développé mon application en react en utilisant la librairie react leaflet et en utilisant une carte open source de openstreetmap.

Pour pouvoir appeler mes fichiers json récupérés en sortie de la chaîne de traitement, j'ai créé un serveur web dans un docker et je les place sur le serveur.

Lorsque mon application react s'exécute, je choisis quel fichier json va être lu. En fonction des id GeoNames **présent** dans mon json, j'exécute une requête sparql vers le Virtuoso afin de récupérer plusieurs informations comme les coordonnées géographiques.

Outre la visualisation cartographique recherchée, les traitements précédents permettent d'enrichir les données à partir des informations obtenues dans Wikipédia, Wikidata et GeoNames

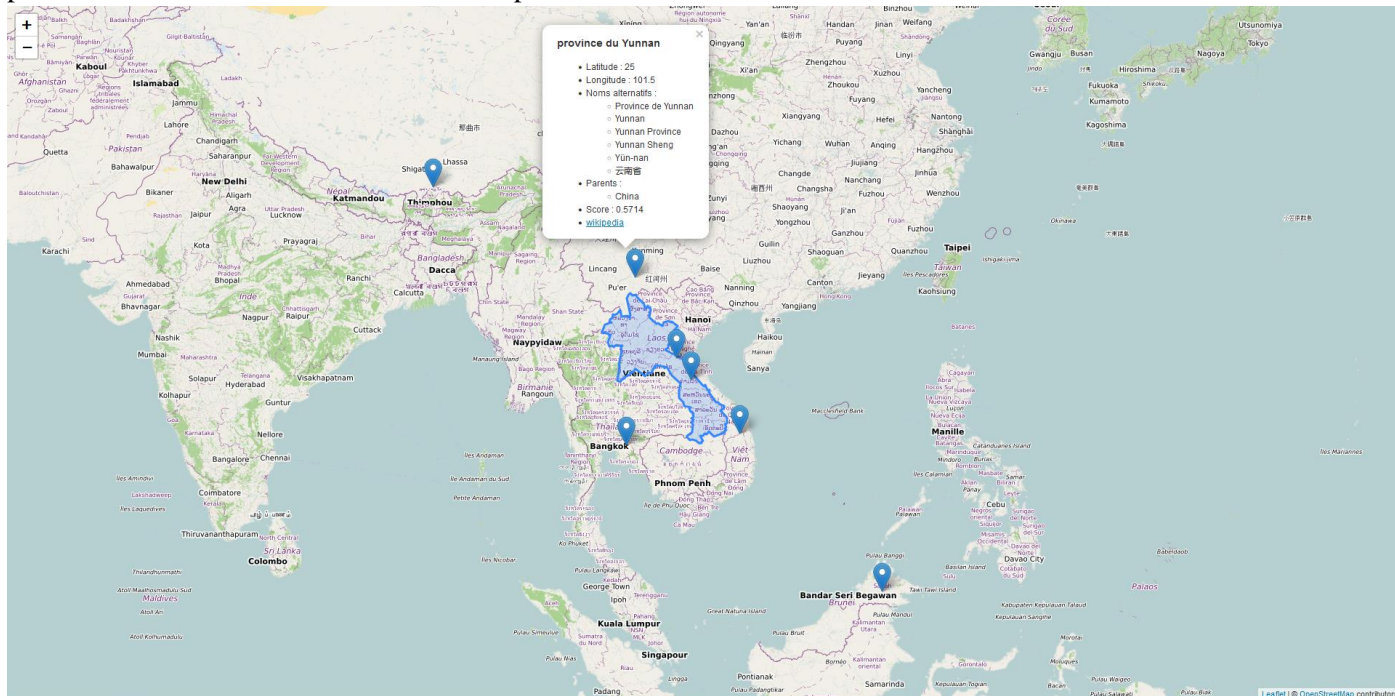


Figure 12 : Visualisation cartographique avec bulle d'information

Dans la bulle d'information, nous pouvons donc trouver un lien vers la page Wikipedia de l'entité géographique pour obtenir plus d'informations.

Dans cet exemple, on peut remarquer que le pays du Laos ne possède pas un marqueur simple mais bien la délimitation de tout le pays (colorié en bleu).

Cela est possible grâce à des fichiers GeoJson qui possèdent un id GeoNames associé au contour d'un pays.

Ces fichiers GeoJSON sont obtenus sur le site de GeoNames.

Pour obtenir ce contour pour des villes, cela n'a pas été possible dans le cadre de ce projet qui utilise prioritairement des données et outils libres et gratuits. Or ce type de données est payant (1500€ par an sur le site de GeoNames).

Conclusion et perspectives

Les objectifs de ce stage étaient de réaliser une chaîne de traitement permettant d'extraire d'un article scientifique les entités nommées géographiques, de les aligner avec un référentiel standard (ici, GeoNames) puis de les visualiser sur une carte.

Ces objectifs font appel à des connaissances en Web sémantique et en traitement automatique des langues. Cela m'a permis de renforcer mes compétences scientifiques et techniques dans ces domaines. J'ai pu aussi mettre en application les bonnes pratiques de développement. Ce dernier point était d'importance pour le projet : le code développé doit pouvoir être réutilisé et intégré à la chaîne de traitement globale.

J'ai été confronté à différents problèmes tout au long de ce stage, à la fois techniques et conceptuels. J'ai appris à les analyser pour trouver des solutions avec l'aide de mes encadrants.

Par ailleurs, j'ai énormément renforcé mes compétences transversales : prise de recul par rapport au sujet, synthèse et rédaction d'un document, relations/communication avec les membres d'un projet.

Si je devais poursuivre mon travail sur ce projet, il y a au moins deux points que j'aimerais approfondir et améliorer :

- Tester ma chaîne de traitement sur un plus grand nombre de textes pour avoir des résultats plus pertinents.
- Trouver une solution (gratuite) pour les fichiers GEOJson afin d'avoir un id GeoNames associé à un contour de ville, pays ou de région.

Dans tous les cas, l'expérience acquise au cours de ce stage me servira grandement dans mes futurs emplois.

Bibliographie

Compréhension générale du sujet

[1] Arnaud, J. Cataloguer, rechercher des cartes. Le référencement géographique en question.
<https://www.cairn.info/revue-documentaliste-sciences-de-l-information-2014-3-page-68.htm#>

[2] Ghislain Auguste Atemezang, Raphaël Troncy. Vers une meilleure interopérabilité des données géographiques françaises sur le Web de données. ».

<https://hal.archives-ouvertes.fr/hal-00716149v1/>

[3] Pascal Cuxac, Alain Collignon, Stéphanie Gegoro, François Parmentier, Des bases de données massives au Web de données : désambiguïsation et alignement d'entités géographiques dans les textes scientifiques.

<https://hal.archives-ouvertes.fr/hal-02307577>

[4] Geographic Ontologies: Survey and Challenges Robert Laurini (Lyon, France) & Okba Kazar (Biskra, Algeria) <https://perso.liris.cnrs.fr/robert.laurini/text/Kazar.pdf>

[5] W3C Geospatial ontologies <https://www.w3.org/2005/Incubator/geo/XGR-geo-ont-20071023/>

Liens utiles

Index géographiques (Gazetteers)

- **GeoNames** : <http://www.GeoNames.org/>

Pour avoir des infos sur l'API : <https://www.programmableweb.com/api/GeoNames>

- **OpenStreetMap & LinkedGeoData**

La base de données associée à OpenStreetMap se trouve sur linkedgeodata.org :

“LinkedGeoData is an effort to add a spatial dimension to the Web of Data / Semantic Web.

LinkedGeoData uses the information collected by the OpenStreetMap project and makes it available as an RDF knowledge base according to the Linked Data principles. It interlinks this data with other knowledge bases in the Linking Open Data initiative.” Visualisation :

<http://browser.linkedgeodata.org/>

- **Getty Thesaurus of Geographic Names**

<https://www.getty.edu/research/tools/vocabularies/tgn/>

- **GeoEthno**, un thésaurus géographique pour l'ethnologie

<http://www.mae.u-paris10.fr/dbtw-wpd/bed/index-lesc.html>

Entity linking avec Wikidata

<http://www.semantic-web-journal.net/system/files/swj2670.pdf>

Documentation et outils

Entity-fishing : <https://nerd.readthedocs.io/en/latest/index.html>

GROBID Named Entity Recognition Documentation : <https://grobid-ner.readthedocs.io/en/latest/> Voir aussi : <https://grobid.s3.amazonaws.com/presentations/29-10-2017.pdf>

GitHub entity-fishing : <https://github.com/kermitt2/entity-fishing>

Service en ligne entity-fishing (API REST) : <http://nerd.huma-num.fr/nerd/>

Client Python: pip install git+<https://github.com/issa-project/entity-fishing-client-python.git>

BLINK: <https://github.com/facebookresearch/BLINK>

Babelfy: <http://babelfy.org>

Docker Virtuoso :

<https://hub.docker.com/r/tenforce/virtuoso/>

React Leaflet :

<https://react-leaflet.js.org/>

GeoJSON :

<https://en.Wikipedia.org/wiki/GeoJSON>

Projet Issa :

<https://issa.cirad.fr/>

Table des figures

Figure 1 : Exemple d'une notice d'un article dans Agritrop	4
Figure 2 : Décomposition du projet issa.....	6
Figure 3 : Résultat Blink avec un texte issu d'agritrop	11
Figure 4: Résultat d'entity-fising avec un texte issu d'Agritrop.....	12
Figure 5 : Exemple de réponse au format json de entity-fishing	13
Figure 6 : Exemple d'une entité nommée géographique non désambiguïsée	14
Figure 7 : Exemple d'une entité nommée désambiguïsée.....	14
Figure 8 : Analyse des résultats pour l'article 598198.....	18
Figure 9 : Analyse des résultats pour l'article 597393.....	19
Figure 10 : Dump de GeoNames avant mise en conformité	21
Figure 11 : Dump de GeoNames après mise en conformité	22
Figure 12 : Visualisation cartographique avec bulle d'information.....	23

Annexes

Config.ini

```
[ENTITY_FISHING]
DOSSIER_ENTREE = C:\Users\Luci\Desktop\Test\txt
DOSSIER_SORTIE = C:\Users\Luci\Desktop\Test\Not_propre
MIN_CONFIDENCE_SCORE = 0.1
API_BASE = http://localhost:5500/service/

[REQUETE]
DOSSIER_ENTREE = C:\Users\Luci\Desktop\Test\Not_propre
DOSSIER_SORTIE = C:\Users\Luci\Desktop\Test\Propre
```

Config.ini qui permet de modifier directement les paramètres importants du code

Entity-fishing.py

```
from nerd import nerd_client
import json
import os
import os.path
import sys
import re
import configparser

config = configparser.ConfigParser()
config.read('config.ini')

#url api
api_base = config.get('ENTITY_FISHING', 'API_BASE')
min_confidence_score = float(config.get('ENTITY_FISHING', 'MIN_CONFIDENCE_SCORE'))
```

1 : Import + config parser

```
#Takes a directory and a file as input and reads the contents of the file
def read(rep, file_name):
    result= None
    with open(os.path.join(rep,file_name),'r',encoding = 'utf-8') as file:
        result = file.read()

    return result
```

2 : Fonction de lecture d'un fichier texte

```
#Depending on the result of the wikidata query and if the type of the entity is LOCATION, keep only the geographical entities and add the Id_geonames
def Traitement_Type_LOCATION(json_,lang):
    result = []

    for data in json_:
        print(' ', data['rawName'])
        if 'wikidataId' in data:
            result.append(data)
        elif data['type'] == 'LOCATION':
            short = entity_fishing_short(data['rawName'],lang)[0]
            if short != None and 'entities' in short:
                data.pop('type', None)
                data.pop('confidence_score', None)
                data['confidence_score'] = short['entities'][0]['confidence_score']
                if 'wikipediaExternalRef' in short['entities'][0]:
                    data['wikipediaExternalRef'] = short['entities'][0]['wikipediaExternalRef']
                if 'wikidataId' in short['entities'][0]:
                    data['wikidataId'] = short['entities'][0]['wikidataId']
            result.append(data)

    return result
```

3 : Traitement des entités nommées avec un type LOCATION

```
#Conforming the json to keep only geographical entities and relevant information. Modifies json display for start and end of words compared to entity-fishing for duplicates
def clean_json(json_):
    result = []
    entities_only = json_[0]['entities']
    for data in entities_only:
        if not any(obj['rawName'] == data['rawName'] for obj in result):
            data.pop('domains', None)
            data['text_location'] = []
            data['text_location'].append({'offsetStart':data['offsetStart'], 'offsetEnd':data['offsetEnd']})
            del data['offsetStart']
            del data['offsetEnd']
            result.append(data)
        else:
            next(obj for obj in result if obj['rawName'] == data['rawName'])['text_location'].append({'offsetStart':data['offsetStart'], 'offsetEnd':data['offsetEnd']})
    return result
```

4 : Suppression des données non importante + modification pour les doublons + sélection des entités nommées géographique (avec id et avec type Location)

```
#Take a directory, a file and write the data contained in the payload
def write(rep, file_name, payload):
    with open(os.path.join(rep, file_name), 'w', encoding = 'utf-8') as file:
        json.dump(payload,file, indent=4, ensure_ascii=False)
```

5 : Fonction pour écrire dans un fichier

```
#allows you to connect to entity-fishing and disambiguate a text
def entity_fishing(t,lang):
    client = nerd_client.NerdClient(api_base)
    test = client.disambiguate_text(t, language=lang)

    return test
```

6 : Fonction qui permet de se connecter à entity-fishing

```
#Allows you to retrieve the language of a text
def entity_fishing_get_language(t):
    client = nerd_client.NerdClient(api_base)
    test = client.get_language(t)

    return test[0]['lang']
```

7 : Permet de récupérer le langage du texte

```
#Allows you to send a single word to disambiguate and link without context, the function needs to know the language of the word to allow disambiguation and linking correctly
def entity_fishing_short(t,lang):
    client = nerd_Client.NerdClient(api_base)
    test = client.disambiguate_query(t, language=lang)

    return test
```

8 : Permet d'envoyer une entité nommée sans contexte à entity-fishing

```
def main(argv):
    dossier_entree = config.get('ENTITY_FISHING', 'DOSSIER_ENTREE')
    dossier_sortie = config.get('ENTITY_FISHING', 'DOSSIER_SORTIE')
    if not os.path.isdir(dossier_entree):
        print("Créer le dossier d'entrée car il n'existe pas encore")
        return
    if not os.path.isdir(dossier_sortie):
        os.mkdir(dossier_sortie)

    for file_name in os.listdir(dossier_entree):
        print(file_name)
        txt = read(dossier_entree, file_name)
        lang = entity_fishing_get_language(txt)
        json_ = entity_fishing(txt,lang)
        json_ = clean_json(json_)
        json_ = Traitement_Type_LOCATION(json_,lang)
        json_ = list(filter(lambda x : x['confidence_score'] >= min_confidence_score, json_))
        print(file_name)
        transform = re.sub(r'.txt$', '.json', file_name)
        json_ = {'lang' : lang, 'entities' : json_}
        write(dossier_sortie, transform, json_)

if __name__ == "__main__":
    main(sys.argv[1:])
```

9 : Le main qui applique toutes les fonctions précédentes + applique le min_confidence_score + vérifie la l'existence des répertoires

Requête.py

```
import json
import os
import os.path
import sys
from string import Template
from SPARQLWrapper import SPARQLWrapper, JSON
from time import sleep
import configparser

config = configparser.ConfigParser()
config.read('config.ini')

#url to make sparql queries on wikidata
url = 'https://query.wikidata.org/sparql'
```

1 : Import + config parser

```
#Takes a directory and a file as input and reads the contents of the file
def read_text(rep, file_name):
    result= None
    with open(os.path.join(rep,file_name),'r',encoding = 'utf-8') as file:
        result = file.read()

    return result
```

2 : Fonction de lecture d'un fichier texte

```
#Takes a directory and a json file as input, reads the contents of the file
def read_json(rep, file_name):
    result= None
    with open(os.path.join(rep,file_name),'r',encoding = 'utf-8') as file:
        result = json.load(file)

    return result
```

3 : Fonction de lecture d'un fichier json

```
#Take a directory, a file and write the data contained in the payload
def write(rep, file_name, payload):
    with open(os.path.join(rep, file_name), 'w',encoding = 'utf-8') as file:
        json.dump(payload,file, indent=4, ensure_ascii=False)
```

4 : Fonction qui permet d'écrire dans un fichier

```
#Makes a query to wikidata. The query is in the txt file "Requete_sparql_wikidata".
def Get_value_id_geonames_wrapper(n, query):
    q = query.substitute(valeur = n)
    sparql = SPARQLWrapper(url)
    sparql.setReturnFormat(JSON)

    sparql.setQuery(q)

    attempt = 0
    MAX_ATTEMPTS = 5

    while attempt <= MAX_ATTEMPTS:
        try:
            ret = sparql.query().convert()
            return ret['results']['bindings'][0]['id_geonames']['value'] if len(ret['results']['bindings']) > 0 else None
        except:
            if attempt < MAX_ATTEMPTS:
                sleep(0.5)
            finally:
                attempt = attempt + 1
    print("        Abbert")
    return None
```

5 : Fonction pour envoyer une requête a Wikidata et récupérer le résultat

```
#Depending on the result of the wikidata query, keep only the geographical entities and add the Id_geonames
def Traitement_wikidata(json_, query):
    result = []

    for data in json_:
        print(' ', data['rawName'])
        if 'wikidataId' in data:
            wikidata_id = Get_value_id_geonames_wrapper(data['wikidataId'], query)
            if wikidata_id != None:
                data['Id_geonames'] = wikidata_id
                result.append(data)

    return result
```

6 : Fonction pour garder uniquement les entités nommées géographiques et rajouter l'id GeoNames obtenue

```
def main(argv):
    dossier_entree = config.get('REQUETE', 'DOSSIER_ENTREE')
    dossier_sortie = config.get('REQUETE', 'DOSSIER_SORTIE')
    query_wikidata= Template(read_text('.', 'Requete_sparql_wikidata.txt'))

    if not os.path.isdir(dossier_entree):
        print("Créer le dossier d'entrée car il n'existe pas encore")
        return
    if not os.path.isdir(dossier_sortie):
        os.mkdir(dossier_sortie)

    for file_name in os.listdir(dossier_entree):
        print(file_name)
        json_ = read_json(dossier_entree, file_name)
        print(' Traitement_wikidata')
        json_['entities'] = Traitement_wikidata(json_['entities'],query_wikidata)
        json_['entities'] = sorted(json_['entities'], key = lambda x : x['confidence_score'], reverse = True)
        print(' Ecriture')
        write(dossier_sortie, file_name, json_)

if __name__ == "__main__":
    main(sys.argv[1:])
```

7 : Le main qui applique toutes les fonctions précédentes + applique le min_confidence_score + vérifie la l'existence des répertoires

Mise_en_confinité_GeoNames_rdf.py

```

import os.path
import os
import re

#Nettoyage d'une ligne
def Clean(l) :
    #Remplace les espaces entre les balises par un saut de ligne, puis les espaces
    l = re.sub(">( *)<", r">\n\<", l)
    #Retire les alternateName et shortName hors fr et en
    l = re.sub("<gn:(alternateName|officialName|shortName) xml:lang=\"(?!fr|en).*>.*</gn:(alternateName|officialName|shortName)>", "", l)
    #Retire les lignes vides
    l = re.sub("\n\s*\n", "\n", l, flags=re.MULTILINE)
    #Retire balise rdf fermante
    l = re.sub("</rdf:RDF>", "", l)
    return l

with open("all-geonames-rdf.txt", "r", encoding='utf-8') as r :
    with open("all-geonames-rdf-propre.txt", "w", True, encoding='utf-8') as w :

        #Première ligne ignorée
        r.readline()

        #Deuxième ligne
        l = r.readline()
        l = Clean(l)
        w.write(l)

        passe = True
        #Reste du fichier
        for line in r:

            #Saute une ligne sur 2 (les urls)
            if(not passe):
                #On retire la balise xml
                line = re.sub("<\?xml version=\"1\.\0\" encoding=\"UTF-8\" standalone=\"no\"?\><rdf:RDF xmlns\cc=\"http://creativecommons.org/licenses/by-sa/4.0/\">", "", line)
                line = Clean(line)
                w.write(line)

            passe = not passe

        w.write("</rdf:RDF>")

```

Mise en conformité du dumps GeoNames

Requête Sparql

```
SELECT ?id_geonames
WHERE
{
    wd:$valeur wdt:P1566 ?id_geonames
}
```

Requête pour récupérer les id GeoNames a partie du Sparql endpoint de Wikidata

```

prefix gn: <http://www.geonames.org/ontology#>
prefix rdfs: <http://www.w3.org/2000/01/rdf-schema#>
prefix wgs84_pos: <http://www.w3.org/2003/01/geo/wgs84_pos#>

select ?name ?officialName ?latitude ?longitude ?altitude ?alternateName ?nameParentFeature ?nameParentCountry
where
{
  <https://sws.geonames.org/1655842/> gn:name ?name;
                                     wgs84_pos:lat ?latitude;
                                     wgs84_pos:long ?longitude;
                                     gn:featureCode ?type.

  OPTIONAL { <https://sws.geonames.org/1655842/> gn:officialName ?officialName}.
  OPTIONAL { <https://sws.geonames.org/1655842/> wgs84_pos:alt ?altitude}.
  OPTIONAL { <https://sws.geonames.org/1655842/> gn:alternateName ?alternateName}.
  OPTIONAL { <https://sws.geonames.org/1655842/> gn:parentCountry ?parentCountry.
            ?parentCountry gn:name ?nameParentCountry}
  OPTIONAL { <https://sws.geonames.org/1655842/> gn:parentFeature ?parentFeature.
            ?parentFeature gn:name ?nameParentFeature}.
}

```

Requête sparql pour extraire les informations du dump GeoNames dans virtuoso

Résumé

Les objectifs de ce stage étaient de réaliser une chaîne de traitement permettant d'extraire d'un article scientifique les entités nommées géographiques, de les aligner avec un référentiel standard (ici, GeoNames) puis de les visualiser sur une carte.

Les tâches réalisées font appel à des connaissances en Web sémantique et en traitement automatique des langues.

Dans ce rapport, nous expliquons comment nous avons conçu et développé une chaîne de traitement permettant d'atteindre ces objectifs : état de l'art des référentiels géographiques, analyse d'outil d'extraction d'entités nommées et visualisation cartographique.

Mots clés : Web sémantique, Web de données, extraction d'entités nommées, traitement automatique des langues, indexation, descripteurs géographiques.

Abstract

The objectives of this internship were to create a processing chain allowing to extract named geographical entities from a scientific article, to align them with a standard reference system (here, GeoNames) and then to visualize them on a map.

The tasks performed call for knowledge in semantic web and automatic language processing.

In this report, we explain how we designed and developed a processing chain to achieve these objectives: state of the art of geographic repositories, analysis of named entity extraction tools and map visualisation.

Keywords: Semantic Web, Web of data, named feature extraction, automatic language processing, indexing, geographic descriptors indexing, geographical descriptors