

Reconnaissance des Entités Nommées spatiales dans un corpus littéraire bruité : Des entités à la carte

Caroline Koudoro-Parfait (1,2,3)

Séminaire STIH

10 Mars 2022



Sens Texte
Informatique
Histoire



caroline.parfait@outlook.fr

(1) OBTA, Sorbonne Université, Paris, France

(2) STIH, Sorbonne Université, Paris France

(3) SCAI, Sorbonne Center for Artificial Intelligence, Paris, France

- 1 Problématique et enjeux
- 2 Constitution du/des corpus
- 3 Mesurer l'impact du bruit de L'OCR sur la REN ?
 - Corrélation entre la qualité de l'image et la sortie OCR ?
 - Comparaisons manuelles
 - Comparaison automatique des résultats des REN
 - Problèmes de mesure : Précision, Rappel, F-score
 - Intersections
 - Distances et similarités : Importance du choix des métriques
- 4 Cartographie
 - Évaluations Cartographiques
 - Projet NER&MAP
- 5 Et après ?

Comment rendre les **systèmes de NER plus robustes** face aux **variations** dans les données qui leur sont soumises par les **utilisateurs** ?

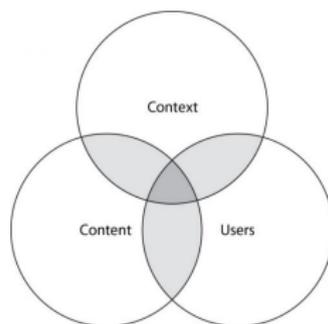


Figure – Peter Morville, *“Three Circles of Information Architecture”*, 2004

- **Utilisateurs** : chercheurs SHS, HN, TAL ... Interdisciplinarité
- **Données** : Corpus ELTeC (Romans Français 19/20ème s.)
- **Variabilités** : diachronie, diatopie, **qualité de la transcription OCR ?**

- Enquêtes **Utilisateurs** quantitative et qualitative et Workshop avec J-B. Tanguy
- **Évaluation des systèmes** Spacy, Stanza, SEM et CasEN
- **Angle Expé.** : "Impact du bruit des transcriptions OCR sur la REN ?"
 - Comment évaluer la qualité des résultats des NER sur des corpus OCR bruités ?
 - Précision, Rappel, F-score ?
 - Approche non supervisée ?

Book	Year	Page nb.
<i>"Le château de Pinon, vol. I "</i> , G. A. Dash, Comtesse	1844	332 p.
<i>"Albert Savarus. Une fille d'Ève. "</i> , Honoré de Balzac	1853	60 p.
<i>"Les trappeurs de l'Arkansas "</i> , Gustave Aimard	1858	450 p.
<i>"Mon village"</i> , Juliette Adam (Lambert)	1860	200 p.
<i>"Le petit chose"</i> , Alphonse Daudet	1868	292 p.
<i>"L'Éducation sentimentale histoire d'un jeune homme"</i> , Gustave Flaubert	1880	520 p.
<i>"Une vie"</i> , Guy de Maupassant	1883	337 p.
<i>"La petite Jeanne"</i> , Zulma Carraud	1884	220 p.
<i>"La Belle rivière"</i> , Gustave Aimard	1894	339 p.
<i>"La nouvelle espérance"</i> , Anna de Noailles	1903	325 p.
<i>"Marie-Claire"</i> , Marguerite Audoux	1925	120 p.

11 ouvrages, 3195 pages

Données

- Littérature, ELTeC¹ Corpus
- Version ELTeC (réf.) and Version OCR
- Comment est constituée la version OCR du corpus?
 - Récupération des images sur le site de Gallica²

1. European Literary Text Collection, <https://www.distant-reading.net/eltec/>

2. <https://gallica.bnf.fr/accueil/fr/content/accueil-fr?mode=desktop>

Constitution du/des Corpus

- Données
 - Littérature, ELTeC¹ Corpus
 - Version ELTeC (réf.) and Version OCR
 - Comment est constituée la version OCR du corpus ?
 - Récupération des images sur le site de Gallica²
- Outils OCR
 - Kraken (Modèle de base)
 - Tesseract (Modèle français et de base)
- Outils de REN
 - Spacy (Modèles français : petit(sm), grand(lg), moyen(md))
 - Stanza (Modèle français)
 - SEM (Modèle WiNER français)
 - CasEN (Modèle de base)

1. European Literary Text Collection, <https://www.distant-reading.net/eltec/>

2. <https://gallica.bnf.fr/accueil/fr/content/accueil-fr?mode=desktop>

Constitution du/des Corpus

- Données
- Littérature, ELTeC¹ Corpus
 - Version ELTeC (réf.) and Version OCR
 - Comment est constituée la version OCR du corpus ?
 - Récupération des images sur le site de Gallica²

- Outils OCR
- Kraken (Modèle de base)
 - Tesseract (Modèle français et de base)

- Outils de REN
- Spacy (Modèles français : petit(sm), grand(lg), moyen(md))
 - Stanza (Modèle français)
 - SEM (Modèle WiNER français)
 - CasEN (Modèle de base)

1. European Literary Text Collection, <https://www.distant-reading.net/eltec/>

2. <https://gallica.bnf.fr/accueil/fr/content/accueil-fr?mode=desktop>

- 1 Problématique et enjeux
- 2 Constitution du/des corpus
- 3 Mesurer l'impact du bruit de L'OCR sur la REN ?
 - Corrélation entre la qualité de l'image et la sortie OCR ?
 - Comparaisons manuelles
 - Comparaison automatique des résultats des REN
 - Problèmes de mesure : Précision, Rappel, F-score
 - Intersections
 - Distances et similarités : Importance du choix des métriques
- 4 Cartographie
 - Évaluations Cartographiques
 - Projet NER&MAP
- 5 Et après ?

Corrélation entre la qualité de l'image et la sortie OCR?



(a) Du bruit dans l'image.



(b) Illustration et légende.



(c) Texte sur deux colonnes.

Figure – Les difficultés de l'OCR³

3. a) "Une vie", Guy de Maupassant. b) "La petite Jeanne", Carraud. c) "Albert Savarus", Balzac

Corrélation entre la qualité de l'image et la sortie OCR ?

Table – Transcriptions OCR d'une image bruitée, première page de "*Une vie*", Guy de Maupassant.

Kraken-Base	Tesseract Français
Texte non transcrit	Texte non transcrit

Corrélation entre la qualité de l'image et la sortie OCR ?

Table – Transcriptions OCR d'une illustration et de sa légende, "La petite Jeanne", Carraud.

Kraken	Tess fr
Ses voisines plumaient leurs oioes quatre fois avant de les ven- LL I II II I I IIF M ii I I II E E g ⁴ Chamnnlhrs de ta mn Mamn- netta ⁵ dre; mais la mere Nan- nette disait que eetaiit une mau- vaise m6thode, paree qu'ainsi la plume ...	Ses voisines plumaient leurs vies quatre fois avant de les ven- Chaumière de la mè1. Nannette ⁶ dre; mais la mère Nannette disait que c'était une mauvaise méthode, parce qu'ainsi l2 plume ...

4. Illustration

5. Légende de l'illustration

6. Légende de l'illustration

Correlation between input quality and output quality?

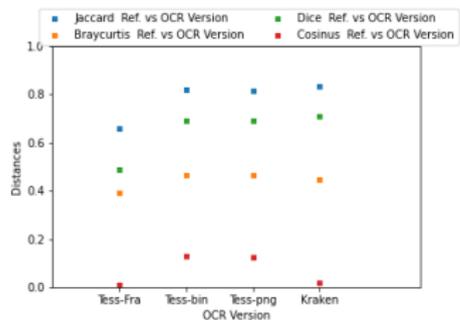
Table – Transcriptions OCR d'un texte mis en page sur deux colonnes, "Albert Savarus", Balzac.

Kraken	Tess fr
Un des quelques salons oh se produisait l'arehe_egue de [] lomene fut l'unique fruit du mariage des Wattoville et des Besangcon sous la Reslauralion, et celui qu'il affectionnait,] de Rupt. etait celui de madame la baronne do Watltev_ile. Un mot] Monsieur de Wotteville passait sn vie dans un riche	Un des quelques salons où se produisait l'archevêque de Besançon sous la Restauration, et celui qu'il affectionnait,[...] Les savans observateurs de la nature sociale ne manqueront pas de remarquer que Phi- ⁷ lomène fut l'unique fruit du mariage des Watteville et des de Rupt ⁸ .

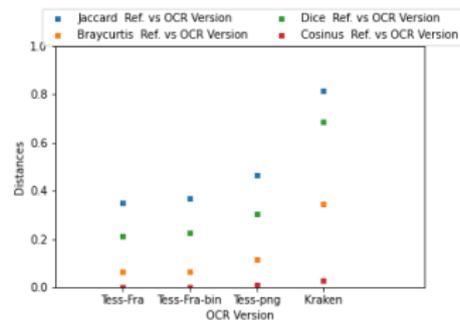
7. Colonne de droite en vert

8. Colonne de gauche en noir

OCR le plus performant ?



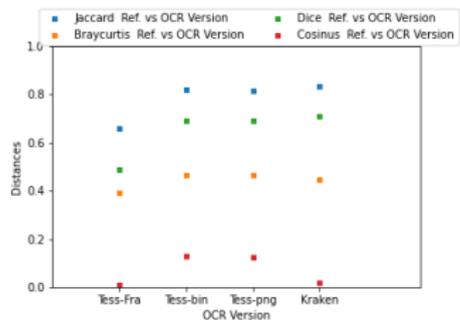
(a) "Albert Savarus", Balzac.



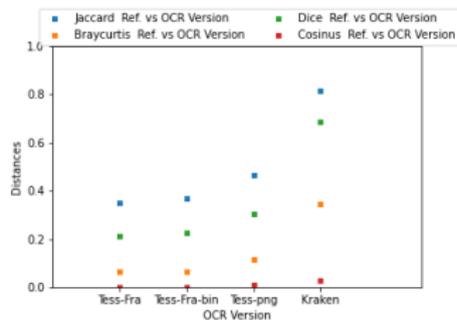
(b) "La petite Jeanne", Carraud.

Figure – Distances entre les versions de Référence et les versions OCR

OCR le plus performant ?



(a) "Albert Savarus", Balzac.



(b) "La petite Jeanne", Carraud.

Figure – Distances entre les versions de Référence et les versions OCR

→ Tesseract semble meilleur que Kraken

Impact des erreurs orthographiques sur la RENS

Version	contexte	spacy_sm	stanza	SEM	CasEN
Réf. ¹	<i>en faïence de Hollande.</i>	Hollande	Hollande	()*	Hollande
Kraken	<i>en faïence _e kollande</i>	()	()	()	()
Tess	<i>en faïence ___ Hollande.. 7 _</i>	Hollande	Hollande	()*	Hollande
Tess fr	<i>en faïence ___ Hollande.</i>	Hollande	Hollande	()*	Hollande
Réf. ²	<i>prendre la diligence de Châlons</i>	Châlons	Châlons	Châlons	()
Kraken	<i>prendre la diligence de Ch_lons</i>	Ch_lons	()	Ch_lons	()
Tess	<i>prendre la diligence de Chalons</i>	Chalons	Chalons	Chalons	()
Tess fr	<i>prendre la diligence de Châlons</i>	Châlons	Châlons	Châlons	()
Réf. ³	<i>on se voit forcé d'aller à Morlincourt.</i>	Morlincourt	Morlincourt	Morlincourt	()*
Kraken	<i>on se voit force d'aller _Mlorlincourt.</i>	Mlorlincourt	Mlorlincourt	Mlorlincourt	()
Tess	<i>on se v01t force (1 al- ' fi.' . . :3 ler £1 Morlincourt.</i>	()	()	Morlincourt	()*
Tess fr	<i>"% ler à Morlincourt.</i>	Morlincourt	Morlincourt	Morlincourt	()*

Table – Substitution ou absence d'un caractère,

¹ "Une vie", Maupassant.

² "L'éducation sentimentale", Flaubert.

³ "Mon village", Adam.

Impact des erreurs orthographiques sur la NER

Version	NER model	Entitiy	# Miss.
Ref.	spacy_lg	Morlincourt : 18	N/A
	STANZA	Morlincourt : 16	N/A
	SEM	Morlincourt : 18	N/A
	CASEN	Morlincourt : 0*	N/A
Kraken	spacy_lg	Morlincourt : 8 Mlorlincourt : 1 Mlorlincourt : 1	8
	STANZA	Morlincourt : 6 Mlorlincourt : 3 Mlorlincourt : 1	6
	SEM	Morlincourt : 8 Mlorlincourt : 3	7
	CASEN	Mlorlincourt : 1	N/A
Tess	spacy_lg	Morlincourt : 11 Morlincourt : 1	6
	STANZA	Morlincourt : 9 Morlincourt : 1	6
	SEM	Morlincourt : 11	7
	CASEN	Morlincourt : 0*	N/A
Tess fr	spacy_lg	Morlincourt : 9 Morlincourt : 1	8
	STANZA	Morlincourt : 7 Morlincourt : 1	8
	SEM	Morlincourt : 11 Morlincourt : 1	6
	CASEN	Morlincourt : 0*	N/A

Table – Variations orthographiques de l'EN "Morlincourt"

- 1 Problématique et enjeux
- 2 Constitution du/des corpus
- 3 Mesurer l'impact du bruit de L'OCR sur la REN ?
 - Corrélation entre la qualité de l'image et la sortie OCR ?
 - Comparaisons manuelles
 - Comparaison automatique des résultats des REN
 - Problèmes de mesure : Précision, Rappel, F-score
 - Intersections
 - Distances et similarités : Importance du choix des métriques
- 4 Cartographie
 - Évaluations Cartographiques
 - Projet NER&MAP

- 5 Et après ?

Précision, Rappel, F-score : problèmes d'alignement

REN Réf.	REN OCR	Verdict 1	Détails	Verdict 2
Oui	Oui	VP	Entité nommée réelle	Vrai VP
Oui	Oui	VP	Erreurs dans les sorties de REN pour la réf. et l'OCR	Faux VP
Non	Non	VN	Aucune entité dans aucune version	Vrai VN
Non	Non	VN	Entité manquante dans les deux versions	Faux VN
Oui	Non	FN	Entité manquante dans la version OCR	Vrai FN
Oui	Non	FN	Erreur d'entité dans la réf.	Faux FN
Non	Oui	FP	Erreur d'entité dans l'OCR	Vrai FP
Non	Oui	FP	Entité manquante dans la réf.	Faux FP (VP)
Non	Oui	FP	Problème de liaison entre entités (<i>entity linking</i>)	Faux FP (VP ?)

Table – Typologie des erreurs de REN.

Problèmes d'entity linking et silence.

Version	Context	Spacy_lg	Stanza	SEM	CasEN
Ref.	[...] la rue Saint-Honoré;	rue Saint-Honoré	rue Saint-Honoré	rue Saint-Honoré	rue Saint-Honoré
Kraken	[...] la rue Saint-Honore;	rue Saint-Honore	rue Saint-Honore	rue Saint-Honore	rue Saint-Honore
Tess	[...] larue Saint-Honoré;	_ Saint-Honoré	()	larue Saint-Honoré	_ Saint-Honoré
Tess fr	[...] la rue Saint-Honoré;	rue Saint-Honoré	rue Saint-Honoré	rue Saint-Honoré	rue Saint-Honoré
Ref.	les États [...] de Guadalajara	Guadalajara	Guadalajara	Guadalajara	Guadalajara
Kraken	les États [...] de Guadalajara	Guadalazara	Guadalazara	Guadalazara	()
Tess	les États [...] de Guadalaxara	Guadalaxara	Guadalaxara	Guadalaxara	()
Tess fr	les États [...] de Guadalaxæw*a	Guadalaæw*a	Guadalaæw*a	()	()

Table – Problèmes d'entity linking : Faux Positif ou Vrai Positif ?¹

Version	spacy_sm	spacy_lg	Stanza	SEM	CasEN
Ref.	Grèce	Grèce bleue (M)	Grèce bleue (M)	Grèce bleue	Grèce
Kraken	Grece	Grece bleue	Grece bleue (M)	()	()
Tess	Grèce (P)	Grèce	Grèce bleue	()	()
Tess fr	Grèce	Grèce	Grèce bleue (M)	Grèce bleue	Grèce

Table – Silence : Vrai négatif ou problème de label ?²

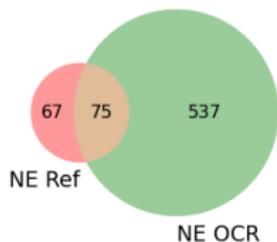
¹ "Le chateau de Pinon", Dash , "Les trappeurs de l'Arkansas", Aimard. ² "La nouvelle espérance", Noailles

Comparaisons Manuelles Vrai Positif/Faux Positif

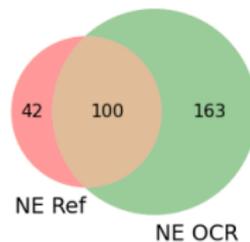
	REF			Kraken		
	SEM	stanza	spacy_lg	SEM	stanza	spacy_lg
TOTAL EN	593	320	273	939	421	946
VP Total	173	179	174	120	137	131
FP Total	415	122	80	811	271	785
AMB Total	5	19	19	8	13	30
Fréq. Total	128	140	142	191	245	612
Répét.	55	40	39	76	48	67
Hapax	73	100	103	115	197	545
VP	67	79	79	49	71	75
FP	58	50	49	135	162	512

Table – "Une vie", Maupassant, 1883.

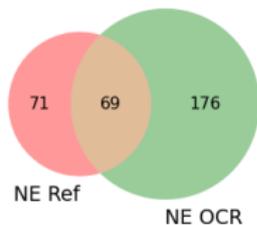
Corrélation entre la qualité de l'OCR et la REN ?



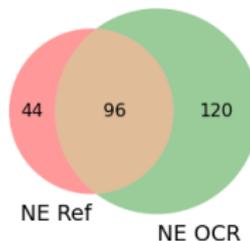
(a) spacy_lg - Kraken



(b) spacy_lg - Tess fr

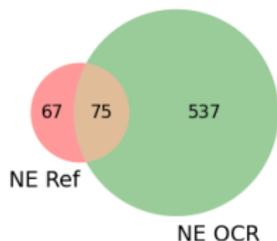


(c) stanza - Kraken

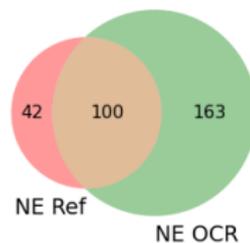


(d) stanza - Tess fr

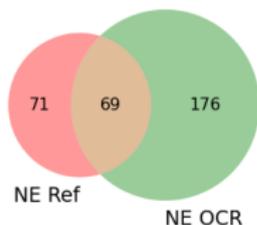
Corrélation entre la qualité de l'OCR et la REN ?



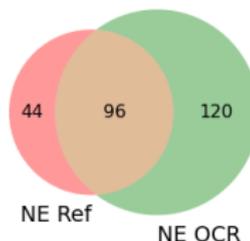
(a) spacy_lg - Kraken



(b) spacy_lg - Tess fr



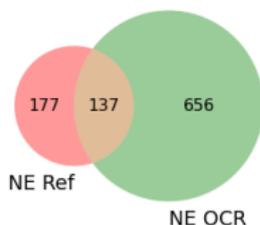
(c) stanza - Kraken



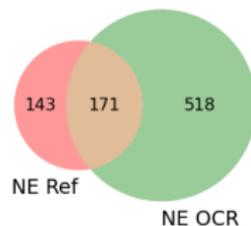
(d) stanza - Tess fr

→ **OCR de meilleure qualité** → **Moins de Faux Négatifs?**⁹

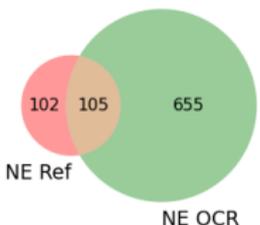
Impact du modèle de NER selon les versions OCR



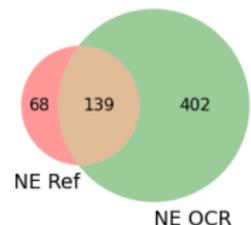
(a) spacy_sm - Kraken



(b) spacy_sm - Tess fr

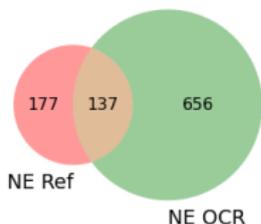


(c) spacy_lg - Kraken

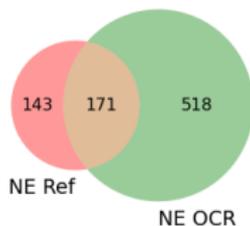


(d) spacy_lg - Tess fr

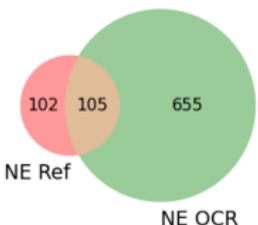
Impact du modèle de NER selon les versions OCR



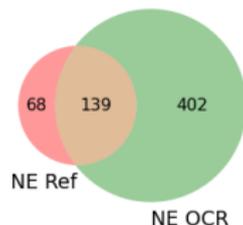
(a) spacy_sm - Kraken



(b) spacy_sm - Tess fr



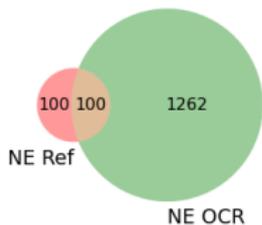
(c) spacy_lg - Kraken



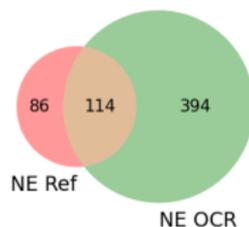
(d) spacy_lg - Tess fr

→ **NER, + de données d'entraînement** → **Moins de FN?** ¹⁰

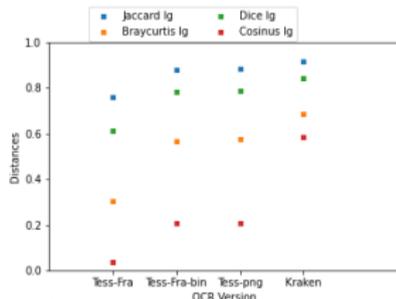
Corrélation entre la qualité de l'OCR et la REN ?



(a) spacy_lg - Kraken



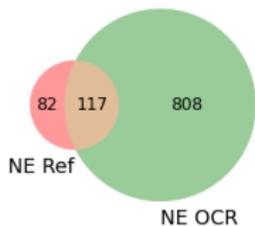
(b) spacy_lg - Tess fr



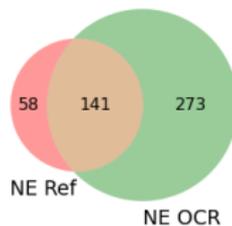
(c) Distances entre les ENS Réf. et les ENS OCR - spacy_lg

11

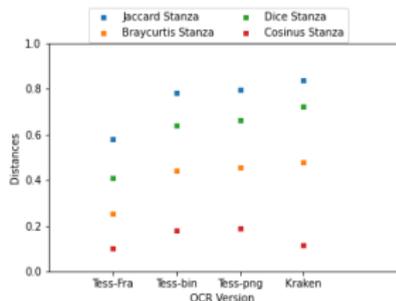
Corrélation entre la qualité de l'OCR et la REN ?



(a) Stanza - Kraken

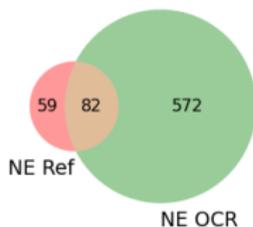


(b) Stanza - Tess fr

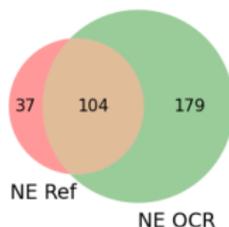


(c) Distances entre les ENS Réf. et les ENS OCR - Stanza
12

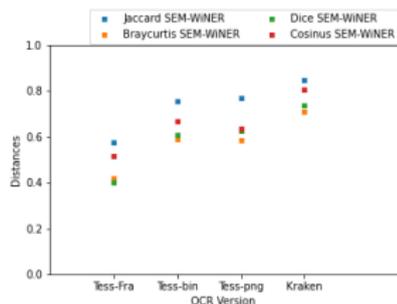
Corrélation entre la qualité de l'OCR et la REN ?



(a) SEM - Kraken



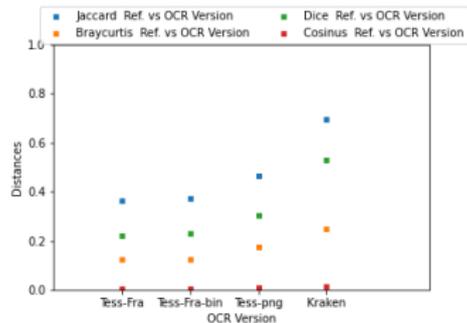
(b) SEM - Tess fr



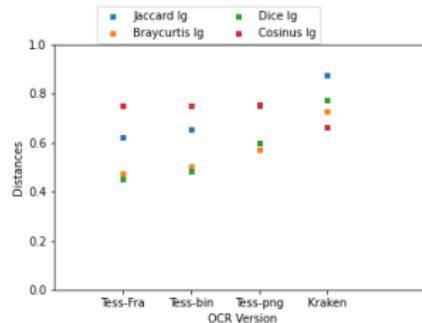
(c) Distances entre les ENS Réf. et les ENS OCR - SEM

→ Comparaison des métriques : à propos du bruit et du silence ? ¹³

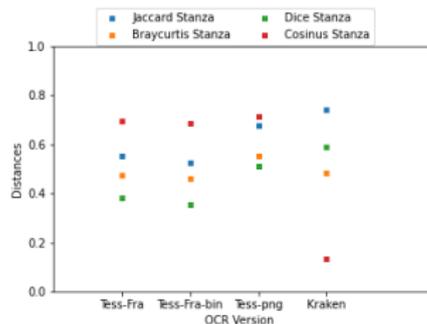
Évaluations des métriques pour les sorties de REN



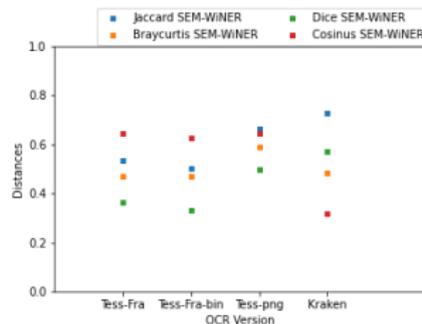
(a) Texte Réf. vs. OCR



(b) `spacy_lg`

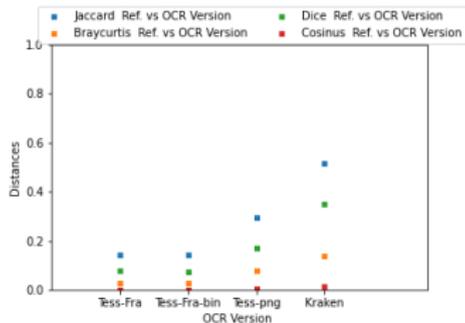


(c) stanza

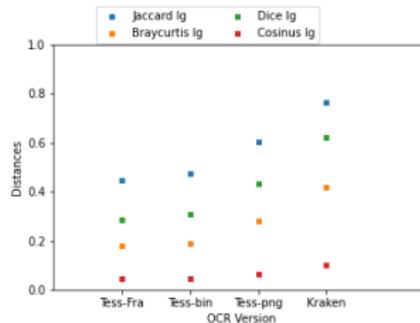


(d) SEM

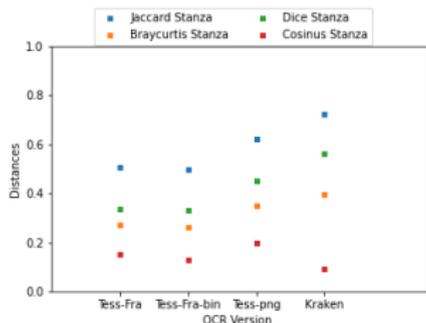
Évaluations des métriques pour les sorties de REN



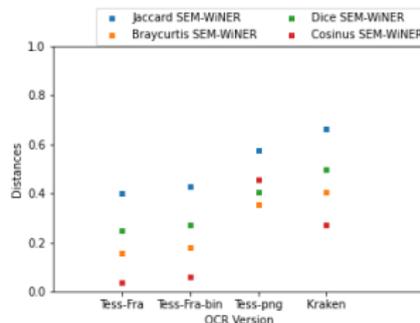
(a) Texte Réf. vs. OCR



(b) spacy_lg

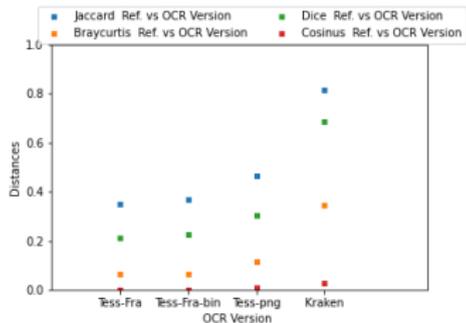


(c) stanza

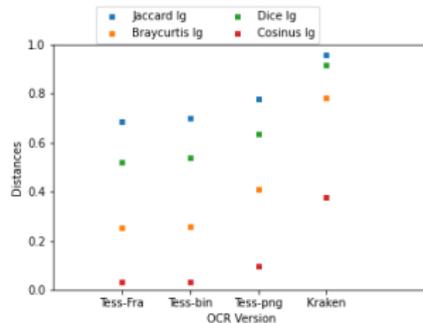


(d) SEM

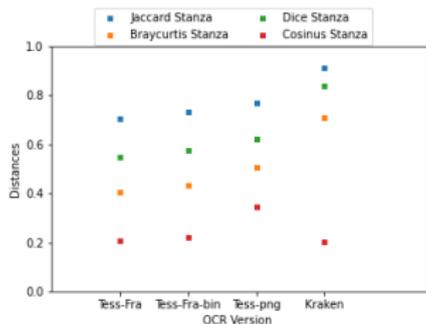
Évaluations des métriques pour les sorties de REN



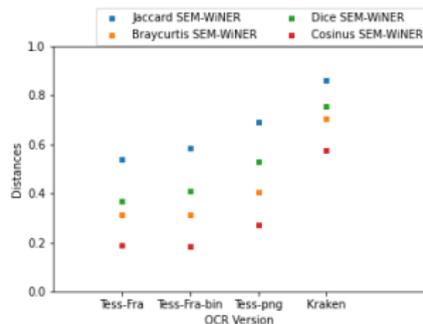
(a) Texte Réf. vs. OCR



(b) spacy_lg



(c) stanza



(d) SEM

- 1 Problématique et enjeux
- 2 Constitution du/des corpus
- 3 Mesurer l'impact du bruit de L'OCR sur la REN ?
 - Corrélation entre la qualité de l'image et la sortie OCR ?
 - Comparaisons manuelles
 - Comparaison automatique des résultats des REN
 - Problèmes de mesure : Précision, Rappel, F-score
 - Intersections
 - Distances et similarités : Importance du choix des métriques
- 4 Cartographie
 - Évaluations Cartographiques
 - Projet NER&MAP

5 Et après ?

Cartographier les résultats OCR basse qualité



(a) Réf. automatique



(b) Réf. corrections sorties REN



(c) Kraken automatique



(d) Kraken corrections sorties REN

Cartographier les résultats OCR bonne qualité



(a) Réf. automatique



(b) Réf. corrections sorties REN



(c) Kraken automatique



(d) Kraken corrections sorties REN

Cartographeur des données imparfaites

← Goderville

id
Goderville

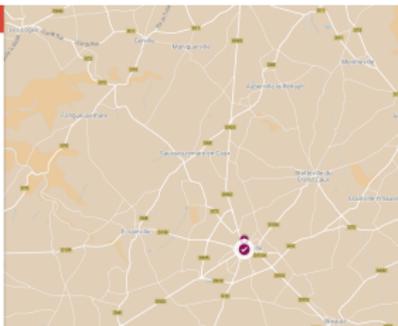
Noms
Goderville ; Goderville ; Goderville

freq_klaxon_spacy_tg
5

latitude
49.64565

longitude
8.36593

GeoNameen_url
<http://www.geonames.org/3013762/goderville.html>



(a) Variations du terme Goderville

← Yport

id
Yport

Noms
Yport ; Yport avec le baron

freq_klaxon_spacy_tg
10

latitude
49.73716

longitude
0.31537

GeoNameen_url
<http://www.geonames.org/2967217/yport.html>



(b) Variations du terme Yport

← Vaucotte

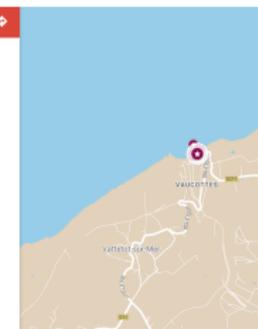
id
Vaucotte

Noms
Vaucotte

freq_klaxon_spacy_tg
1

latitude
49.738276

longitude
0.2918113



(c) Variations du terme Vaucotte

← Vaucotte

id
Vaucotte

Noms
Vaucotte

freq_klaxon_spacy_tg
1

latitude
49.7382373

longitude
0.2908388



(d) Variations du terme Vaucotte

Évaluer visuellement les résultats

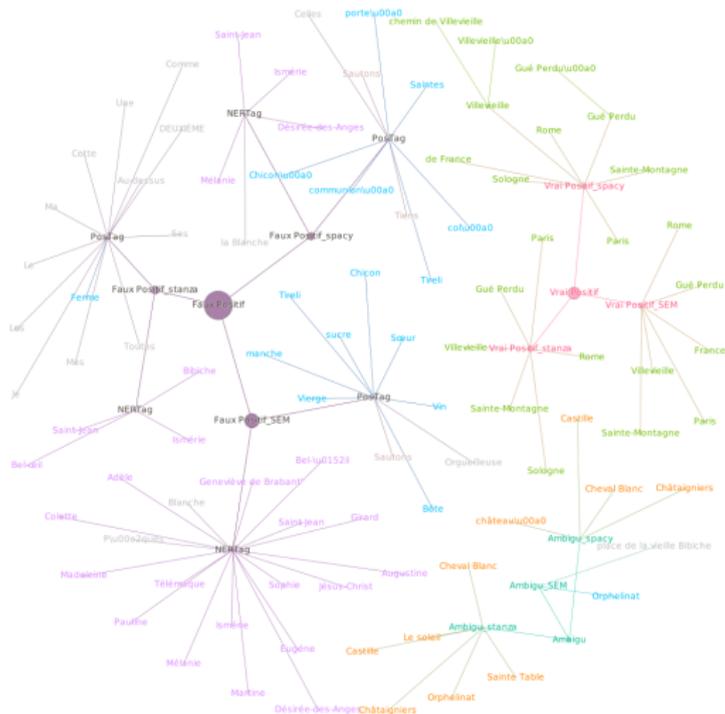
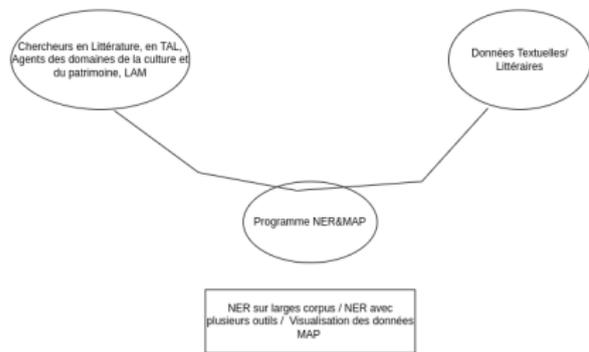
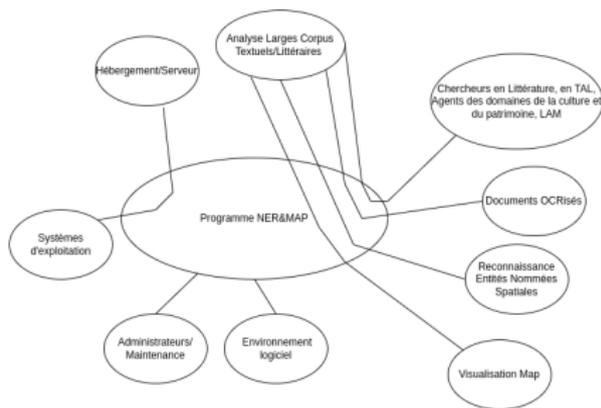


Figure – Visualisation des Vrais Positifs et Faux Positifs Spacy, Stanza et SEM



(a) Bête à corne - Spécifications



(b) Pieuvre - Spécifications

Cas Usage NER&MAP

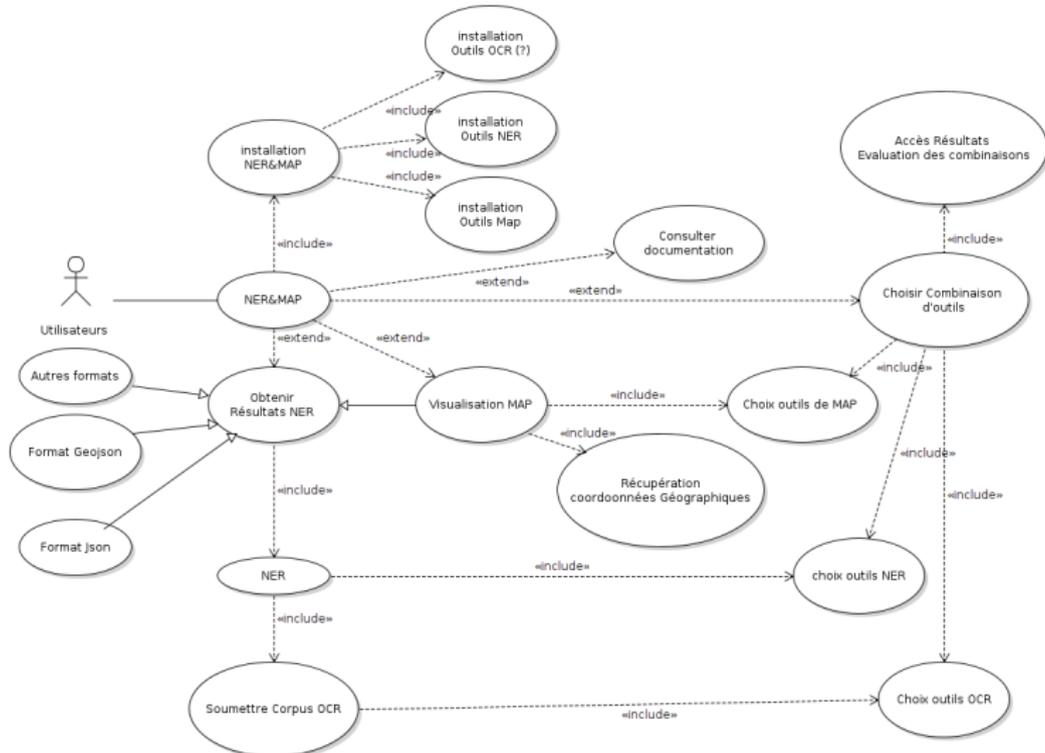


Figure – Cas d'usage

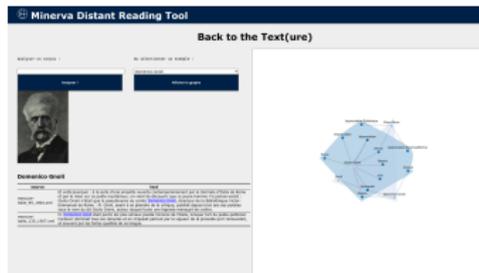
Modèle Application



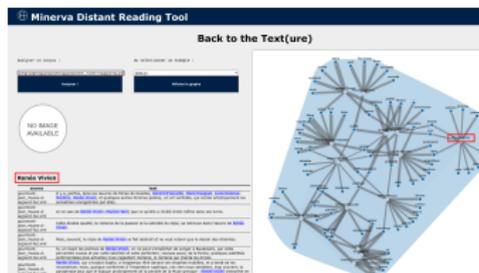
(a) Utiliser l'url d'un fichier xml



(b) Utiliser les données pré-enregistrées



(a) Affichage Graph et métadonnées



(b) Cliquer sur un noeud pour afficher les métadonnées

Figure – Minerva, Y.Dupont

Wireframe NER&MAP

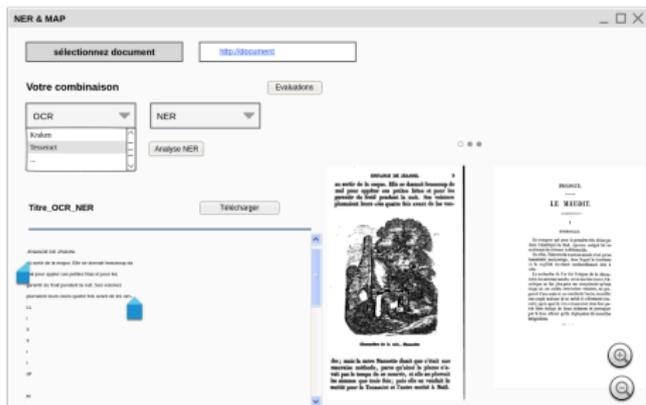


(a) Accueil

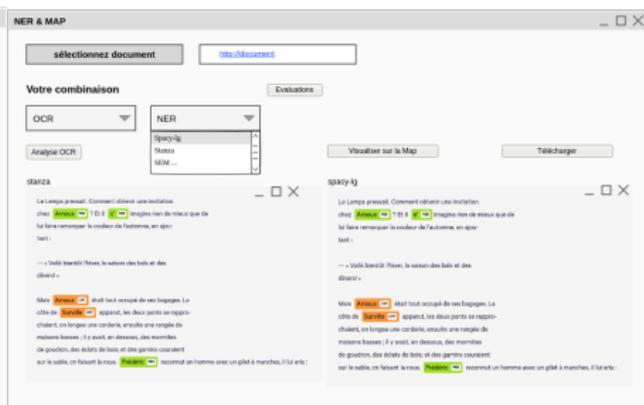


(b) Exemples évaluations

Wireframe NER&MAP



(a) Paramètre OCR



(b) Paramètre NER

Wireframe NER&MAP

NER & MAP

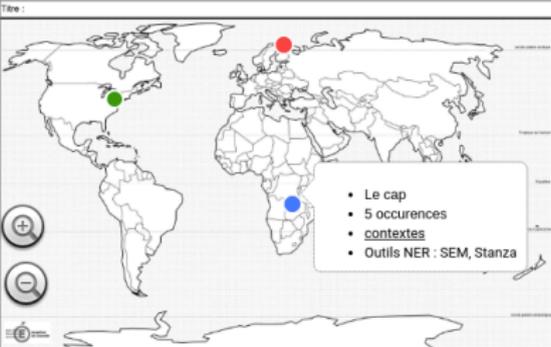
sélectionnez document

Votre combinaison

OCR NER

Visualiser sur la Map

Titre :



- Le cap
- 5 occurrences
- contextes
- Outils NER : SEM, Stanza

● 3 outils de NER
● 2 outils de NER
● 1 outils de NER

Titre_OCR_NER

Entité LOC	Contexte	coord. geo
Paris	Paris, la capitale Française	48.8619, 2.2938
Parisl	Il se rendait à Parisl	<input type="text" value="Proposition de coord"/>

Figure – Visualisation Cartes

	GGMymAP	Leaflet	Qgis
Avantages	Facile d'utilisation, customisation, intégration web facile	Customisation, nombreux tuto ; plusieurs formats d'entrée (Geojson, csv) ; intégration web	Plusieurs couches de cartes
Inconvénients	Les données sont sur Google	R et Rstudio windows	Coût d'entrée élevé

Table – Premiers retours sur différents outils de Visualisation

- Toutes les erreurs d'OCR ne se valent pas
 - fautes d'orthographe \neq mots collés ensemble
- Impact du bruit OCR sur les sorties de NER
 - Variations orthographiques d'un toponyme : comment les lier ?
 - comment les cartographier ?
- Le choix des métriques n'est pas trivial
 - Distance Cosinus : sous-évalue les différences(?)
 - Distance de Jaccard : sur-évalue les différences(?)
 - Proximité distances Cosinus et Jaccard : Problème ?
- La représentation cartographique permet de visualiser :
 - EN récupérée par +ieurs outils est potentiellement un VP ?
 - Une géographie littéraire du 19ème siècle
- Les humains peuvent gérer le bruit, mais que faire s'il y a du silence ?

- Résolution des problèmes d'entity linking (diachronie, formes fautives des mots)
- Améliorer la Géolocalisation automatique
- Annotation et entraînement d'un modèle sur des données bruitées
→ Pour quel résultat ? Meilleur, moins bon, équivalent ?