

# Reconnaissance d'entités nommées sur des sorties OCR bruitées : des pistes pour la désambiguïisation morphologique automatique.

**Caroline Koudoro-Parfait (1,2,3)**, Gaël Lejeune (2), Richy Buth (2)

TALN | RECITAL 2022 - Atelier TALN & HN 2022

27 Juin 2022



Sens Texte  
Informatique  
Histoire



caroline.parfait@sorbonne-universite.fr  
gael.lejeune@sorbonne-universite.fr  
buth\_richy@hotmail.com

(1) OBTr, Sorbonne Université, Paris, France

(2) STIH, Sorbonne Université, Paris France

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

- 1 Problématique(s) et enjeux
- 2 Constitution du/des corpus
- 3 Influence de la qualité de l'image sur la sortie OCR ?
- 4 Observation des entités contaminées
- 5 Évaluations automatiques et visualisations des résultats ?
  - Intersections
  - Nerval : Précision, rappel, f1
  - Importance des métriques : Levenshtein vs. Cosinus
- 6 Et après ?

- "Impact du bruit des transcriptions OCR sur la REN ?"
  - Comment évaluer la qualité des résultats des REN sur des corpus OCR bruités ?
    - Précision, Rappel, F-score ?
    - Approche non supervisée ?
- Il est possible de récupérer des EN sur des textes bruités :
  - Comment associer les EN contaminées à leur forme standard ?

- Données
- Corpus ELTeC<sup>1</sup> :
    - Littérature Française 19ème siècle,
    - 11 ouvrages, 3195 pages

---

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

## Données

- Corpus ELTeC<sup>1</sup> :
  - Littérature Française 19ème siècle,
  - 11 ouvrages, 3195 pages

## Outils OCR

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

---

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

- Données
- Corpus ELTeC<sup>1</sup> :
    - Littérature Française 19ème siècle,
    - 11 ouvrages, 3195 pages
- Outils OCR
- Kraken (Modèle de base)
  - Tesseract (Modèle français et de base)
- Outils de REN
- Spacy (small, medium et large)
  - Stanza
  - SEM (WiNER)
  - CasEN

---

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

# Constitution du/des Corpus

- Données
- Corpus ELTeC<sup>1</sup> :
    - Littérature Française 19ème siècle,
    - 11 ouvrages, 3195 pages

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

- Outils de REN
- Spacy (small, medium et large)
  - Stanza
  - SEM (WiNER)
  - CasEN

- Alignement & Éval.
- Nerval : évaluation de la REN sur du texte bruité

---

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

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

# Influence de la qualité de l'image sur la sortie OCR?



(a) Du bruit dans l'image.

ENTRÉE AU CHAUMON.  
4  
se met de la coupe. Elle se dressa brusquement  
mal pour appeler ses petites filles et pour les  
garantir du froid pendant la nuit. Ses vêtements  
plissaient sous son quatre bras ouvert de son sein.



Chaumotte de la nuit. Maman

des; mais la mère Maumotte disait que c'était une  
mauvaise méthode, parce qu'autre la place n'a-  
vait pas le temps de se sécher, et elle se plissait  
les épaules que trois fois; puis elle se vendait la  
moitié pour le Toussaint et l'autre moitié à Noël.

(b) Illustration et légende.



(c) Texte sur deux colonnes.  
Figure – Les difficultés de l'OCR<sup>2</sup>

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

# Influence de la qualité de l'image sur la sortie OCR?

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- <sup>3</sup> lomène fut l'unique fruit du mariage des Watteville et des de Rupt <sup>4</sup> .

3. Colonne de droite en vert
4. Colonne de gauche en noir

# 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	()	()

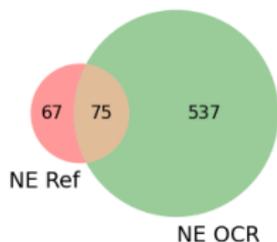
Table – Problèmes d'entity linking : Faux Positif ou Vrai Positif ?<sup>1</sup>

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

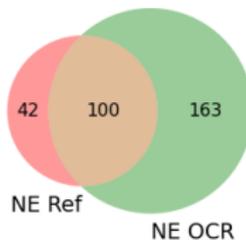
Table – Silence : Vrai négatif ou problème de label ?<sup>2</sup>

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

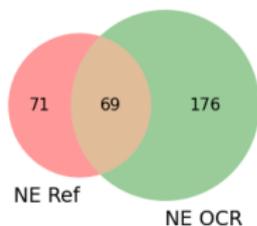
# Corrélation entre la qualité de l'OCR et la REN ?



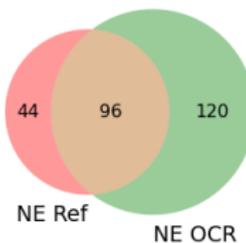
(a) spacy\_lg - Kraken



(b) spacy\_lg - Tess fr

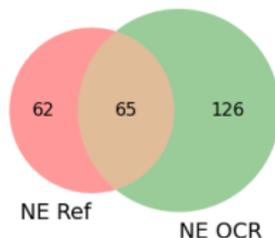


(c) stanza - Kraken

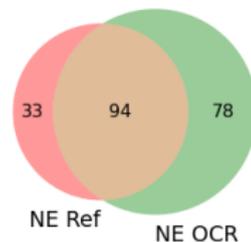


(d) stanza - Tess fr

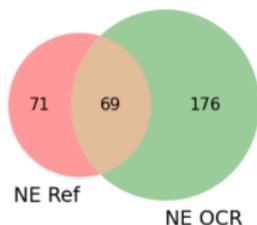
# Corrélation entre la qualité de l'OCR et la REN ?



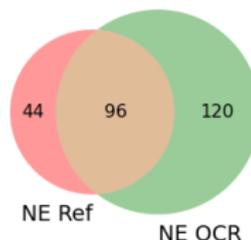
(a) SEM - Kraken



(b) SEM - Tess fr



(c) stanza - Kraken



(d) stanza - Tess fr

→ **OCR de meilleure qualité** → **Moins de Faux Négatifs?**<sup>5</sup>

5. "Une vie", Maupassant, 1883.

# Évaluation et alignement avec Nerval

Version	#Entités		Évaluation par NERVAL			
	Version OCR	Référence	Intersection	Précision	Rappel	$F_1$ mesure
Kraken	1391	965	576	0.414	0.597	0.489
Tess fr	980	965	713	0.728	0.739	<b>0.733</b>
Tess	1090	965	608	0.558	0.630	0.592

Table – Comparaison des résultats de la reconnaissance d'entités nommées avec `spacy_lg` sur différentes versions de "*Le petit chose*", Daudet, 1868, après alignement avec NERVAL

Version	#Entités		Évaluation par NERVAL			
	Version OCR	Référence	Intersection	Précision	Rappel	$F_1$ mesure
Kraken	364	100	41	0.113	0.410	0.177
Tess fr	289	100	49	0.170	0.490	<b>0.252</b>
Tess	433	100	47	0.109	0.470	0.176

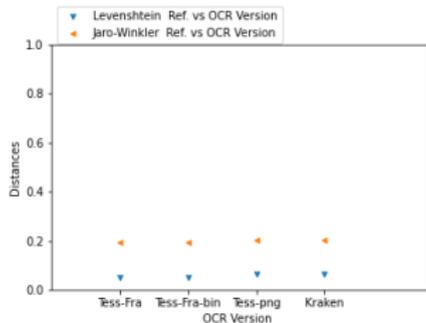
Table – Résultats de la reconnaissance d'entités nommées par système sur différentes versions de "*Marie-Claire*", Audoux après alignement avec NERVAL

# Liage des Entités contaminées avec NERVAL

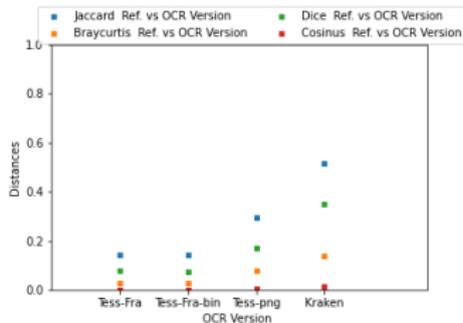
Versions	Entités Réf. <sup>I</sup>	Entités Align	Entités Réf. <sup>II</sup>	Entités Align
Kraken Tess fr Tess	l'Amérique	— merique <b>dé</b> Amérique -Amérique	Paris	<b>eares</b> <b>-aran</b> <b>aran</b>
Kraken Tess fr Tess	rio Gila	Rio Gila Bio Gz- rio Gila	Laon	Laou Laon Laon

**Table** – Alignement des entités (spacy\_1g) de la référence et des versions OCR avec NERVAL, <sup>I</sup>"Les trappeurs de l'Arkansas", Aimard, 1858 et <sup>II</sup>"Le château de Pinon, vol. I.", Dash, 1844

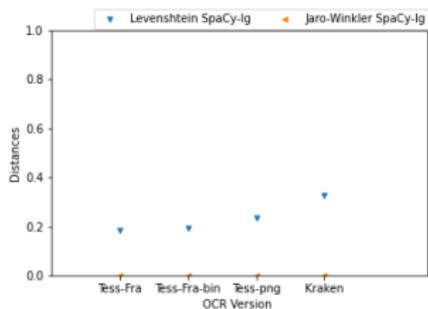
# Importance des métriques : Levenshtein vs. Cosinus



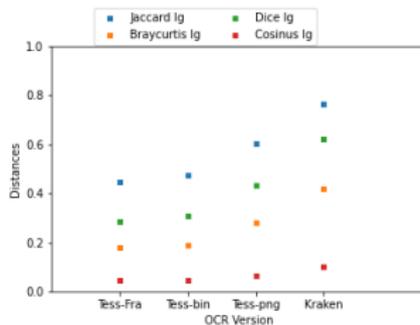
(a) Réf. vs OCR



(b) Réf. vs OCR



(c) spacy\_lg



(d) spacy\_lg

Figure – Distances (gauche) de Levenshtein et Jaro-Winkler entre les Entités Nommées "LOC" Réf. et OCR, "Le petit chose", Daudet, 1868.

# Des pistes pour le liage des Entités contaminées

Versions	Entités Réf.	Entités Align	Jaccard <sup>6</sup>	Cosinus <sup>6</sup>
Kraken	Morlincourt	Mlorlincourt	0.1428	0.0715
		Mlorlincourtl	0.1818	0.1210
Tess fr	Morlincourt	Morlinco'urt	0.1818	0.0762
		Morlin	<b>0.4761</b>	<b>0.2788</b>
Kraken	Saint-Brunelle	Brunclle	<b>0.5925</b>	<b>0.3244</b>
		Brunelle	0.4583	0.2012
Tess fr	Saint-Brunelle	Saint—Brunelle	0.2222	0.0909
		Saint—anelle	0.4642	0.2183

Table – Résultats de la récupération automatique des entités contaminées récupérées par spacy\_1g sur différentes versions de "Mon village", Adam, 1860.

## 6. Bigrammes de caractères

- 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 ?
- Le choix des métriques n'est pas trivial
  - Distance Cosinus : évalue la différence de quantité entre deux groupes(?)
  - Distance de Jaccard : sur-évalue les différences(?)
  - Levenshtein vs. Cosinus : résultats proches, cos. moindre coût de calcul
- Les humains peuvent gérer le bruit, mais que faire s'il y a du silence ?

- Résolution des problèmes de liage (diachronie, formes fautives des mots)
- Désambiguïsation automatique : liage avec des bases de données ou embeddings ?