

# Exploration de la collection Française d'ELTeC avec l'outil de cartographie textuelle Epiméthée

---

Caroline Koudoro-Parfait  
caroline.parfait@sorbonne-universite.fr

20 Mars 2025

Observatoire des Textes des Idées et des Corpus - Obtic,  
Sorbonne Center for Artificial Intelligence - SCAI,  
Sens Textes Informatique Histoire - STIH EA 4509, Sorbonne Université

# Plan de la présentation

Enjeux, usages et verrous à surmonter

Impact des contaminations OCR sur la REN

Stratégies pour aider les utilisateur·ice·s

Explorer l'espace littéraire européen

Et après ?

## Enjeux, usages et verrous à surmonter

---

## Utilisateurs et interdisciplinarité :

- Sciences Humaines et Sociales (SHS) et Lettres ;
- Traitement automatique des langues (TAL) ;

## ➔ Acquisition des corpus

→ Reconnaissance optique de caractère (OCR) ;

## Variabilité dans le contexte d'usage

→ **qualité de la transcription OCR ;**

- Comment adapter les outils de TAL aux pratiques des chercheurs en Littérature ?
- Quelle chaîne de traitement pour assister les chercheurs en Lettres ?
- Quelles stratégies pour dépasser les contaminations OCR ?

## Données : Textes Littéraires

\* European Literary Text Collection (ELTeC)<sup>1</sup>  
[Schöch et al., 2021] :

- 22 langues,  $\approx$  100 romans par langue ;
- période : 1840 à 1920 ;

---

1. <https://www.distant-reading.net/eltec/>

## ✂ Corpus constitués

Corpus	Books	Pages	Words	# Named Entities : LOC		# Évaluations		
				spaCy_lg	flair	qti.	qli	Cluster
<i>small-ELTeC-fra</i>	11	3 195	829 604	5 765	4 814	✓	✓	✓
<i>small-ELTeC-eng</i>	9	5 281	2 063 246	5 551	8 867	✓	✓	✓
<i>small-ELTeC-por</i>	4	1 795	421 915	7 590	N/A	✓	✓	✗
Total	24	10 271	3 314 765	13 906	13 681	-	-	-

**Table 1** – Statistiques sur les corpus dans la version de référence

# Méthodes explorées et leurs performances

Rechercher des solutions appropriées pour rendre la REN utilisable sur des données bruitées :

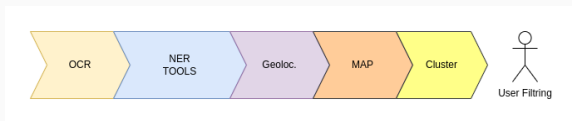
- **Désambiguïisation**, en utilisant des **métriques de similarité** et des **clusters** ✓ [Koudoro-Parfait et al., 2022]
- **Combinaison**, plusieurs modèles de REN pour filtrer les résultats ✓ [Petkovic et al., 2025]
- **Correction automatique**, problèmes de sur-correction ✗ [Koudoro-Parfait et al., 2024]
- **Linking**, Lier des entités contaminées à leur forme standard ✗
- **Cartes** Une vue d'ensemble des entités ✓ [Koudoro-Parfait and Lejeune, 2024]

Performances : ✓ très convaincantes ; ✓ convaincante ; ✓ peu convaincantes.



# Épiméthée, agir avant de réfléchir !

- ☞ OCR → REN – sans correction de la sortie OCR ;
- ☞ Intégration stratégies de filtrages auto. des FP de REN :
  - Combiner des systèmes de REN ;
  - Clustering ;
  - Géolocalisation ;



**Figure 2** – Chaîne de traitement Épiméthée,  
<https://github.com/These-SCAI2023/EPIMETHEE>

- ☞ Interface : filtre manuel par l'utilisateur-ice ;
- ☞ Récupération (CSV), réemploi et analyse littéraire.

# Impact des contaminations OCR sur la REN

---

# Difficultés rencontrées par les modèles d'OCR



(a) Illustration + légende



(b) Texte en colonnes



(c) Texte en filigrane



(d) Décoration, Capitalisation

Figure 3 – a) G. de Maupassant, *Une vie*, 1883. b) Inconnu, *Adélaïde de Mariendal*, drame en cinq actes, 1783. c) Z. Carraud, *La petite Jeanne*, 1884. (d) H. de Balzac, *Albert Savarus*, 1853.

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

Kraken 3.0	Tess. fr 0.3.6
Ses voisines plumaient leurs oies quatre fois avant de les ven- LL I I I I I I I F M ii I I I E E g Chamnlnhrs de ta mn Mamnnetta dre ; mais la mere Nannette disait que eétait une mauvaise m6thode, paree qu'ainsi la plume ...	Ses voisines plumaient leurs vies quatre fois avant de les ven- Chaumière de la mè1. Nannette dre ; mais la mère Nannette disait que c'était une mauvaise méthode, parce qu'ainsi l2 plume ...

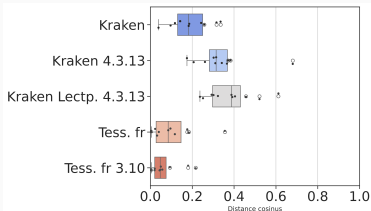
**Table 2** – Transcriptions\* OCR d'une illustration et de sa légende. ● = illustration, ● = légende, ● = contaminations orthographiques.

\* Z. Carraud, *La petite Jeanne*, 1884.

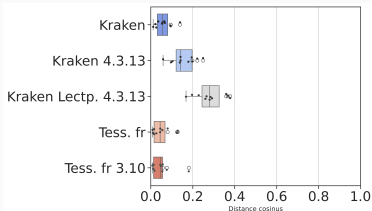
## Entités contaminées → Faux Positif ou Vrai Positif contaminé

	Context	spaCy_lg	Flair
Ref	<i>There are some fields near Manchester</i>	Manchester	Manchester
Krak	<i>_ are some fields near <b>ance</b>hester [...]</i>	ancehester	()
Tess. en	N/A	N/A	N/A
Ref	<i>toute la noblesse de Picardie</i>	Picardie	Picardie
Krak	<i><b>toutela</b> noblesse de <b>V</b>icardie</i>	<b>V</b> icardie	<b>V</b> icardie
Tess. fr	<i>toute la noblesse de Picardie</i>	Picardie	Picardie

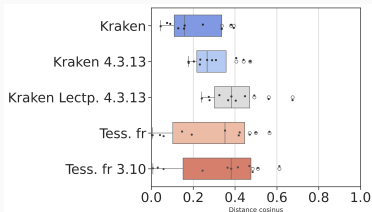
# Meilleure configuration selon cosinus



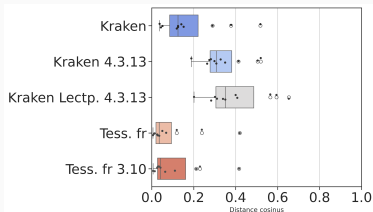
(a) `spacy_lg`, cosinus



(b) `Flair`, cosinus.



(c) `stanza`, cosinus



(d) `CamemBert`, cosinus

**Figure 4** – Distances cosinus pour la REN sur le corpus *small-ELTeC-fra* (global).

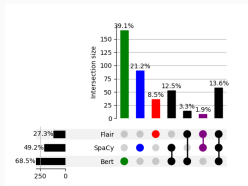
- Toutes les erreurs d'OCR ne se valent pas ;
- Meilleures configurations : Tesseract – spaCy\_lg ou Tesseract – flair ;

# Stratégies pour aider les utilisateur·ice·s

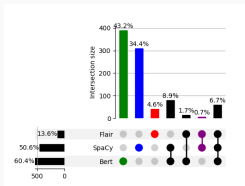
---



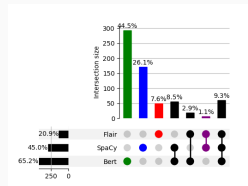
# Combinaison de $\neq$ systèmes de REN



(a) Référence



(b) Kraken



(c) Tess\_fr

**Figure 5** – Combination of three NER tools on high quality OCR (Tess.fr : CER = 0.03, WER = 0.05 ; Kraken : CER = 0.05, WER = 0.18). Each column represents the percentage of NEs found by each combination of tool(s) (exclusive). Whilst each row gives the individual, and inclusive, coverage of each subset. Daudet, 1868.

# Cluster : des pistes pour le liage des entités contaminées ?



**Figure 6** – Cluster cosinus bigramme de caractères sur les sorties Tesseract, spaCy\_lg. A. Daudet, *Le petit chose*, 1868.

# Cluster : des pistes pour le liage des entités contaminées ?

→ Clusters intéressants, le centroïde et un VP.

Version	Centroid	Cluster members
Réf. <sup>a</sup>	Montparnasse	Montparnasse, boulevard Montparnasse, théâtre Montparnasse, Montmartre, rue Bonaparte, Mont-, Saumon, Gymnase
Kraken	Montparnasse	Montparnasse, boulevard Montparnasse, theatre Montparnasse, Gymnase, Debarrassez, WWt3, rs5, ytP
Kraken <sup>b</sup>	Goderville	Gdoderville, Gloderville, Goderville, Barville, Fourville, OD0

→ Mais parfois le centroïde peut être un FP.

Version	Centroid	Cluster members
Réf. <sup>a</sup>	PION	Lyon, Odéon, Rio, PION
Kraken <sup>a</sup>	Fougeroux	Luxembourg, Perou, Broum, Fougeroux, MY, Vaudoux, lesFougeroux

- erreur Clustering,
- erreur Clustering + interférence OCR,
- erreur Clustering + bruit REN,
- EN LOC.

<sup>a</sup> modèle REN spaCy\_1g, "Le petit chose", Daudet, 1868.

<sup>b</sup> modèle REN spaCy\_1g, "Une vie", G. de Maupassant, 1883

	<i>small-ELTEC-fra</i>	<i>small-ELTEC-eng</i>
Ref.	329	93
Kraken	828	356
Tess.	1035	116
Total	2192	565

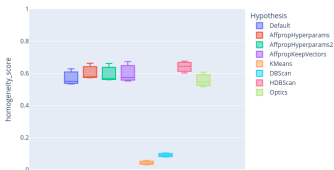
**Table 3** – Statistics on the annotated subset

Les algo. évalués + CountVectorizer, n-gram(min, max)

- Affinity Propagation :
  - Default (2,2);
  - Hyperparams (2,4);
  - Hyperparams2 (3,4);
  - KeepVectors (3,4);
- DBScan (2,4);
- HDBScan (2,4);
- Optics (2,4);

# Évaluation quantitative de $\neq$ algo. de clustering

homogeneity\_score by hypothesis for en



(a) homogeneity score, en.

homogeneity\_score by hypothesis for fr



(b) homogeneity score, fr.

completeness\_score by hypothesis for en



(c) completeness score, en.

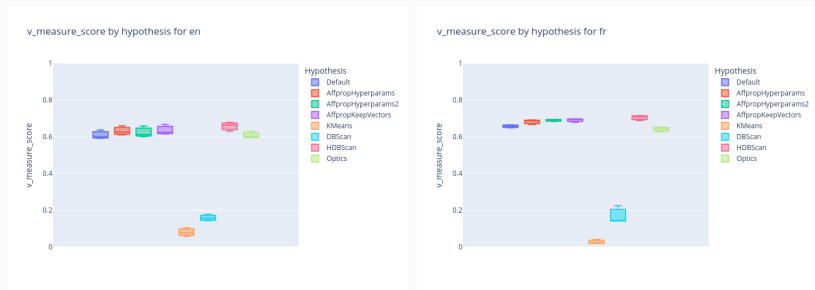
completeness\_score by hypothesis for fr



(d) completeness score, fr.

Figure 7 – Evaluation of clustering algorithms for *small*-ELTeC-eng and

# Évaluation quantitative de $\neq$ algo. de clustering

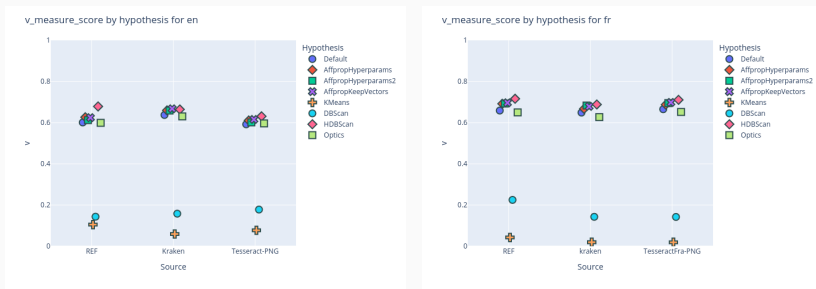


(a) v-measure, en.

(b) v-measure, fr.

**Figure 8** – Evaluation of clustering algorithms for *small*-ELTeC-eng and *small*-ELTeC-fra.

# Évaluation quantitative de $\neq$ algo. de clustering



**Figure 9** – Evaluation of various clustering algorithms on different versions of text (reference and OCR) using v-measure, for the *small-ELTeC-eng* and *small-ELTeC-fra*.



# Explorer l'espace littéraire européen

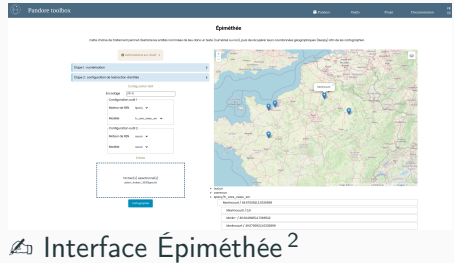
---

- ⇒ Modèles d'OCR : Tesseract Français, Anglais, Portugais ;
- ⇒ Modèles de REN par langues :

	fr	en	pt
spaCy	✓	✓	✓
Flair	✓	✓	✗

# Épiméthée : un système de bout-en-bout

- 👉 Conception du projet ;
- 👉 Conduite du projet ;
- 👉 Gestion équipe ;



📄 Interface Épiméthée<sup>2</sup>

---

2. <https://github.com/These-SCAI2023/EPIMETHEE>

# Épiméthée : un système de bout-en-bout

Bati	Batignolles
Bati	Bati
Batignolles	Batignolles
Cuba	Cuba
TOUR	Batignolles
Annou	Cuba
Autour	Luxembourg
Luxembourg	TOUR
TOUR	Touraine
Touraine	Turc
Turc	

Batignolles
Bati
Batignolles
Cuba
Luxembourg
Touraine
TOUR
Annou
Autour
TOUR
Turc

(a) Output from kraken

(b) Output from kraken man. corr.

**Figure 10** – Results of the Épiméthée pipeline with the spaCy\_lg and flair models, on kraken version. Before and after user filtering *man. corr.* = manually corrected by the user ; ● centroid error compared to candidates, ● candidate error, ● manual filtering.

Le cas "Arthurville", *Capitaine Cap*, A. Allais.

- spaCy\_lg et flair = OK
- Épiméthée absent  $\Rightarrow$  qui est coupable :
  - Géoloc. ?  $\Rightarrow$  Actuel Municipalité de Saint Raphaël (Québec).
  - Clustering ? *Aff. Prop.* écarte des entités.

$\rightarrow$  Combien d'autres cas ?

$\rightarrow$  Quelles stratégies d'évaluations ?

Cartes obtenues à partir de la chaîne de traitement Épipiméthée et filtrage sur les romans de :

- Hector Malot, "Sans famille", 1878 ;
- Pierre Loti, "Mon frère Yves", 1883 ;
- Catulle Mendès, "Luscignole", 1892 ;
- Alphonse Allais, "Capitaine Cap", 1902.

[https://drive.google.com/drive/folders/  
1h3sUycwQKexjMuC01bLMJdTEsvJRiIDD?usp=drive\\_link](https://drive.google.com/drive/folders/1h3sUycwQKexjMuC01bLMJdTEsvJRiIDD?usp=drive_link)

**Et après ?**

---

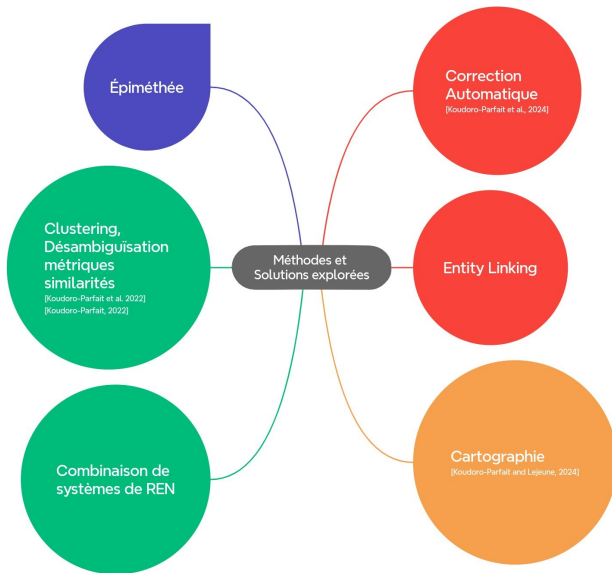
- ➔ Éval. impact erreurs OCR sur REN → **tâche non triviale** :
- ➔ Épiméthée → *end-to-end* pour dépasser les contaminations OCR ;
- ➔ Évaluation d'algo. de clustering et paramétrages → amélioration des **clusters** ;
- ➔ Contribuer à une géographie littéraire du XIX<sup>e</sup> siècle : Quel orient chez les Orientalistes ? ;



- ➔ Un texte *propre* **ne garantit pas** de meilleurs résultats de REN :
  - Correction auto. et **sur-correction** ;
  - **Erreurs** des outils de REN même sur textes *propre* ;
  - Importance choix **métrique** éval. ;
- ➔ Assistance utilisateur-ice-s :
  - **Combinaison** systèmes, EN extraite par +ieurs outils → VP ;
  - **Clustering** : rapprocher des formes contaminées d'une EN ;

- Intégrer : observations de la **fréquence** documentaire et **longueur** ;
- Liage avec des bases de données de toponymes anciens ;
- Améliorer la Géolocalisation auto. ;

# Merci de votre attention !





Koudoro-Parfait, C. and Lejeune, G. (2024).

**Reconnaissance des Entités Nommées spatiales sur un corpus littéraire bruité : des entités à la carte.**



*In Séminaire des sources aux Systèmes d'Information Géographique.*




Koudoro-Parfait, C., Lejeune, G., and Buth, R. (2022).

**Reconnaissance d'entités nommées sur des sorties ocr bruitées : des pistes pour la désambiguïsation morphologique automatique.**

*In TAL-HN @ TALN(Traitement Automatique des Langues Naturelles) 2022.*

-  Koudoro-Parfait, C., Petkovic, L., and Roe, G. (2024).  
**Analyse multilingue de l'impact de la correction automatique de la roc sur la reconnaissance d'entités nommées spatiales dans des corpus littéraires.**  
*Revue TAL.*  
À paraître.
-  Petkovic, L., Koudoro-Parfait, C., Desmarest, M.-S., and Lejeune, G. (2025).  
**Quelle solution pour améliorer les performances de la reconnaissance d'entités nommées sur des données bruitées, corriger l'entrée ou filtrer la sortie ?**  
*Corpus, (26).*

-  Schöch, C., Patras, R., Erjavec, T., and Santos, D. (2021). **Creating the european literary text collection (eltec) : Challenges and perspectives.**  
*Modern Languages Open*, 0(1) :25.

# Corpus ELTeC - European Literary Text Collection

Ouvrage	Auteur	Année	Pages	Mots	spaCy_lg	stanza
<i>Mon village</i>	J. Adam	1860	200	20938	213	152
<i>La Belle rivière</i>	G. Aimard	1894	339	137392	1004	959
<i>Les trappeurs de l'Arkansas</i>	G. Aimard	1858	450	91119	646	606
<i>Marie-Claire</i>	M. Audoux	1925	120	35780	101	108
<i>Albert Savarus. Une fille d'Ève</i>	H. de Balzac	1853	60	79924	682	684
<i>La petite Jeanne</i>	Z. Carraud	1884	220	53212	316	95
<i>Le château de Pinon, vol. I</i>	G. A. Dash	1844	332	44246	271	311
<i>Le petit chose</i>	A. Daudet	1868	292	86482	744	580
<i>L'Éducation sentimentale</i>	G. Flaubert	1880	520	150494	1304	1098
<i>Une vie</i>	G. de Maupassant	1883	337	75745	302	312
<i>La nouvelle espérance</i>	A. de Noailles	1903	325	54272	182	236

**Table 5** – *small-ELTeC-fra*, 11 ouvrages, 3195 pages

# Corpus ELTeC - European Literary Text Collection

Ouvrage	Auteur	Année	Pages	mots	spaCy_lg	stanza
<i>Home influence</i>	G. Aguillar	1847	628	171342	205	244
<i>Auriol</i>	W. H. Ainsworth	1844	246	46388	82	55
<i>Wuthering Heights</i>	E. Brontë	1847	764	94986	140	132
<i>Coningsby</i>	B. Disraeli	1844	983	101778	634	543
<i>Mary Barton</i>	E. Gaskell	1848	423	161568	290	281
<i>The Mysteries of London</i>	G. Reynolds	1844	840	810167	2019	2312
<i>Modern Flirtations vol.1</i>	C. Sinclair	1841	386	189057	502	248
<i>Vanity Fair</i>	W. M. Thackeray	1848	624	298568	1492	1164
<i>The Life and Adventures of M. Armstrong</i>	F. Trollope	1840	387	189392	187	207

**Table 6** – *small-ELTeC-eng*, 9 ouvrages, 5281 pages



Ouvrage	Auteur	Année	Pages	Mots	spaCy_lg	stanza
<i>Quattro Novelas</i>	A. Castro Osorio	1908	272	50766	353	N/A
<i>A illustre casa de Ramires</i>	E. de Queirós	1900	543	107441	3881	N/A
<i>O crime do padre Amoro</i>	E. de Queirós	1875	620	141700	2362	N/A
<i>Uma familia ingleza</i>	J. Diniz	1875	360	122008	994	N/A

**Table 7** – *small-ELTeC-por*, 4 ouvrages, 1795 pages

# La Très grande Bibliothèque (TGB)

Ouvrage	Auteur	Année	Pages	Tokens	spaCy_lg	stanza
<i>La princesse Pallianci</i>	C. L. Bazan-court	1852	340	36 423	304	263
<i>Meryem, scènes de la vie algérienne. Marcel.</i>	C. Perrier [Bentégeat]	1863	360	85 077	662	512
<i>Wilmina, ou L'enfant des Apennins</i>	L. G. de Caude-berg	1820	242	36218	353	188
<i>Les fourmis du parc de Ver-sailles raisonnant ensemble dans leurs fourmilières</i>	C. Lambert	1803	72	10 173	57	47
<i>Œuvres complètes de Pierre Loti</i>	P. Loti	1893-1911	588	133 129	2040	2136
<i>La confession d'un enfant du siècle / Alfred de Musset ; avec un portrait... par Eugène Lami...</i>	A. de Musset	1879	494	92 140	578	269
<i>Le Parnasse envahi, petit poème allégorique au sujet du sacre de S. M. Charles X.</i>	E. Rullier	1825	71	10 261	165	38
<i>La Comtesse de Rudolstadt</i>	G. Sand	1861	340	102 423	618	505
<i>Diégarias, drame en 5 actes et en vers</i>	V. Séjour	1844	38	18 603	970	293
<i>Le département de l'Oise : Compiègne et Marat, frag-ment historique</i>	A. Sorel	1865	19	6 277	108	105

**Table 8** – *small-TGB-RevuesCorpus*, 10 ouvrages, 2564 pages.

# La Très grande Bibliothèque (TGB)

Ouvrage	Auteur	Année	Pages	Tokens	spaCy_lg	stanza
<i>L'Alsace et la Lorraine</i>	L. Longret	1873	2	357	13	12
<i>La Grèce libre</i>	A. Bignan	1821	20	1 027	35	19
<i>Poésies diverses</i>	Inconnu	1745	10	1 502	32	11
<i>Les dernières Étrivières [...]</i>	B. Bonafoux	1877	22	2 320	29	20
<i>M. de L'Espinasse [...]</i>	D. L. Baric	1851	20	3 058	102	91
<i>Adélaïde de Mariendal, drame en cinq actes</i>	Inconnu	1783	100	15 344	276	217
<i>Œuvres du seigneur de Brantôme. Tome 14</i>	P. de Bourdelle Sgr de Brantôme	1779	255	49 084	844	507
<i>Souvenirs d'un vieux mélomane</i>	A. Pontmartin	1879	350	61 872	659	598
<i>La lyre des petits enfants</i>	A. Cordier	1857	357	62 639	646	447

**Table 9** – *small-TGB-RevueTAL*, 9 ouvrages, 1136 pages.

# Les hapax : indice de la contamination de l'OCR

→ Loi de Zipf, transcription qualité ↗, CER : Kraken = 0.0886, avec Tess. fr = 0.0496

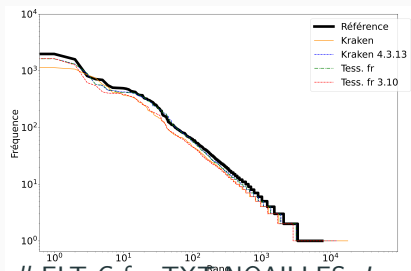


Figure 11 – *small-ELTeC-fra* TXT NOAILLES, *La nouvelle esperance*

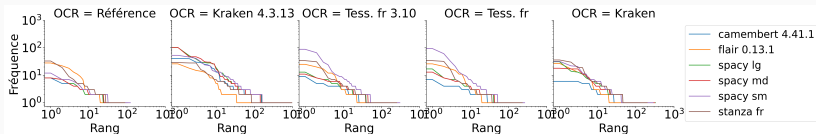


Figure 12 – *small-ELTeC-fra* REN NOAILLES, *La nouvelle esperance*