

Réutiliser les corpus numériques constitués par OCR : bruit et silence dans les data papers

Caroline Koudoro-Parfait (1,2,3), Gaël Lejeune (2)
Jean-Baptiste Tanguy (1,2)

DHNord 2021

17 Novembre 2021



Sens Texte
Informatique
Histoire



caroline.parfait@sorbonne-universite.fr,
jean-baptiste.tanguy@sorbonne-universite.fr
gael.lejeune@sorbonne-universite.fr
(1) OB TIC, Sorbonne Université, Paris, France
(2) STIH, Sorbonne Université, Paris France
(3) SCAI, Sorbonne Center for Artificial Intelligence, Paris, France

Récupération de données à partir de transcriptions automatiques

→ Choix de l'OCR? Le data Paper propose :

➤ Tableau d'observation manuelle

➤ Evaluation automatique des sorties OCR

Kraken	Tess fr
__ P_TITB JLANN_ IL DRVOIR. PHEIIikl PAIII_ EM_ANCE DE JEANME, _a mbre Nannette.	__ PETITE JEANNE LE DEVOIR. _ PREMIÈRE PARTIE. EN- FANCE DE JEANNE. La mère Nannette.

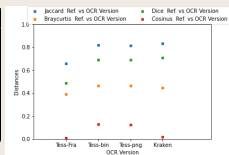
Table – 2 sorties OCR sur le texte * ELTEC.

Kraken 17
8 .2 S De pouuoir vous lauuer la pance En despit de toute la France, [...] Vos bonnets a la Polonnoife lfe Que l'on rajufte a la Francoife ;l, ile 4 A Par _s, le dernier du mois , je Six io nrs deuant le iour des loils _ Sesif a2 4 . FIN. ↵

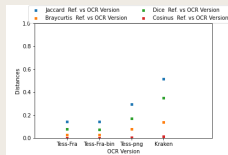
Table – 1 sortie OCR sur le texte **
ANTONOMAZ.

Bruit et silence OCR :

- Rouge : substitutions de caractères.
- Bleu : transcription d'un graphisme.
- " _ " : caractère manquant.

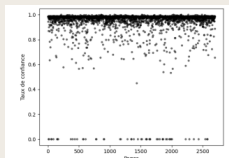


(a) "Albert Savarus", Balzac, 1853.

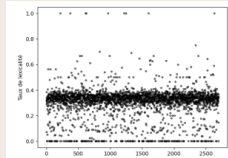


(b) "Le petit chose",
Daudet, 1868.

Figure – Mesurer la qualité de l'OCR, métriques :
Jaccard, Cosinus, Dice et Braycurtis
(1= Différent, 0= Similaire).



(a) Taux de confiance.



(b) Taux de Lexicalité.

Figure – Mesurer la qualité de l'OCR pour tout
le corpus Antonomaz (1= Qualité correcte).

REN sur des corpus bruités : Observation manuelle des données

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.

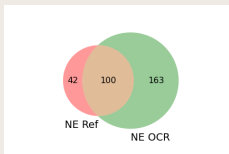
Version	Context	spacy_1g	stanza
Ref.	[...] la rue Saint-Honoré;	rue Saint-Honoré	rue Saint-Honoré
Kraken	[...] la rue Saint-Honore;	rue Saint-Honore	rue Saint-Honore
Tess	[...] larue Saint-Honoré;	_ Saint-Honoré	()
Tess fr	[...] la rue Saint-Honoré;	rue Saint-Honoré	rue Saint-Honoré
Ref.	les États [...] de Guadalajara	Guadalajara	Guadalajara
Kraken	les États [...] de Guadalazara	Guadalazara	Guadalazara
Tess	les États [...] de Guadalaxara	Guadalaxara	Guadalaxara
Tess fr	les États [...] de Guadalaæw*a	Guadalaæw*a	Guadalaæw*a
Ref.	La Rochelle.[...] Thouloufe.	La Rochelle, () *	La Rochelle, Thouloufe
Kraken 17	La Rocdhelle .[...] Thouloufe.	Rocdhelle , () *	La Rocdhelle , Mer. Thouloufe
Ref.	[...]Viennes en Auftriche.	Viennes, Auftriche	Viennes, Auftriche
Kraken 17	[...]Viennes en Auftriche.	Viennes, Auftriche	Viennes, Auftriche

Table – Problèmes d'entity linking et silence.

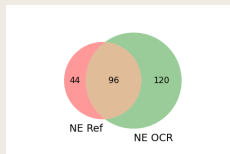
*Silence dû au s long non utilisé dans la graphie actuelle? cf. exemple suivant.

Évaluations des sorties de REN dans le data paper

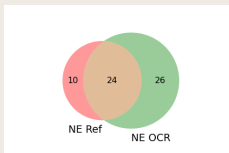
➤ Intersections :



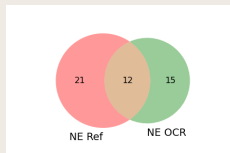
(a) spacy_1g - Tess fr



(b) stanza - Tess fr



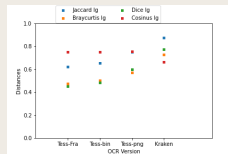
(c) spacy_1g - Kraken 17



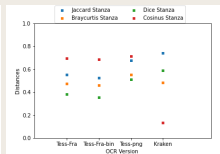
(d) stanza - Kraken 17

- NE OCR : Problèmes d'entity linking et bruit corrigéables.
- NE REF : Vrai positif (potentiel silence dans la référence?).
- Intersection : Vrai positif et erreurs Réf./OCR de l'outil REN.

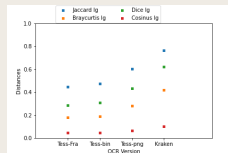
➤ Distances et Similarités :



(a) spacy_1g



(b) stanza



(c) spacy_1g



(d) stanza

- Distance Cosinus : sous-évalue les différences(?)
- Distance de Jaccard : sur-évalue les différences(?)
- Proximité distances Cosinus et Jaccard : Problème?

a b "Une vie", Maupassant, 1883.

c d "Le petit chose", Daudet, 1868.