

Correction automatique des interférences OCR dans la reconnaissance d'entités nommées spatiales : réel gain ou perte de l'information ?

Caroline Koudoro-Parfait^{1,2,3}, Ljudmila Petkovic¹, Gaël Lejeune²

Journée d'étude « Bruit de fond ou valeur ajoutée ? »

Grenoble, 28 avril 2023



SensTexte
Informatique
Histoire



- (1) caroline.parfait@sorbonne-universite.fr
- (2) ljudmila.petkovic@sorbonne-universite.fr
- (3) gael.lejeune@sorbonne-universite.fr

- (1) OBtic, Sorbonne Université, Paris, France
- (2) STIH, Sorbonne Université, Paris, France
- (3) SCAI, Sorbonne Center for Artificial Intelligence, Paris, France

- 1 Contexte et enjeux de la REN sur des données transcrites par OCR
 - Influence de la qualité de l'image sur la sortie OCR ?
 - Observation manuelle des EN contaminées issues des OCR corrigés et non corrigés
- 2 Évaluation automatique de l'impact de la correction OCR sur la REN
 - Intersections : EN réf. vs. EN OCR + EN réf. vs. EN OCR corrigés
 - Distances et similarités : Importance du choix des métriques
 - Nerval : outil d'alignement des EN contaminées & rappel, précision, f1
- 3 Et après ?

Question Faut-il corriger les erreurs d'OCR ou tenter de rendre les **systèmes de REN plus robustes** face aux **variations** dans les données qui leur sont soumises par les **utilisateurs** ?

- Angles Expé.**
- « Impact du bruit des transcriptions OCR sur la REN ? »
Comment évaluer la qualité des résultats de REN sur des corpus OCR bruités ?
 - « Impact de la correction des transcriptions OCR sur la REN ? »

Données Litt. Française 19^e siècle,

- Corpus ELTeC (Réf.)¹ : 10 ouvrages, 3195 pages.
- Très grande bibliothèque (TGB)² : 10 ouvrages, 1157 pages.

OCR Kraken, Tesseract (Modèle français et de base)

Modèles REN fra. Spacy, Stanza, SEM, CasEN

- Corr. auto**
- JamsPELL³
 - Un autre ?

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

2. TGB, <http://obvil.lip6.fr/tgb/>

3. <https://github.com/bakwc/JamSpell>

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

Table – Transcriptions OCR d'une page de texte théâtral, avec décoration, "*Drame en cinq actes*", Inconnu.

Kraken	Tess fr
A D E L A I D E v reee= e,E=eeee=ees,g s c E N E I I I. a E LE BARON, LA B ARONNE. r (Pendant l_ derniere Scene la Baronne a temoigne te grand aattement.) L s B A R O N. _E vois, Madame, que vos regrets egalent les miens. [...] L A B A R O N _ _ . Quand on perd fon bien, on ofe tout.e L B B A R O N. Ils favent que je fuis honnete homme, & _ _ auroient grand tort...	6 'ADÉLAIDE SCENE III. LE BARON, LA BARONNE. (Pendant la derniere Scene la Baronne a témoigné un grand abattement,) Le BARON. \$ E vois, Madame, que vos regrets égalent les miens.[...] LA BARONNE. Quand on perd fon bien, on ofe tout, LE Baron. Ils favent que je fuis honnête homme, & ils auroient grand tort... -°

→ Difficultés de l'OCR : Bleu, Décorations ; Vert, Capitalisation ; Rouge, bruits divers. Bruits = ajout/transformation/suppression

Alignement manuel d'EN contaminées, OCR non corrigés

Version	Context	spaCy_lg	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 Etats [...] de Guadalazara	Guadalazara	Guadalazara
Tess	les Etats [...] de Guadalaxara	Guadalaxara	Guadalaxara
Tess fr	les États [...] de Guadalaæw*a	Guadalaæw*a	Guadalaæw*a
Ref.**	[...] rue de Sèvres;	rue de Sèvres	rue de Sèvres
Kraken	[...] ruc de Sbvrcs	Sbvrcs	Sbvrcs
Tess	[...] rue de Sévres	Sévres	rue de Sévres
Tess fr	[...] rue de Sèvres;	rue de Sèvres	rue de Sèvres
Ref.**	le faubourg Saint-Germain	faubourg Saint-Germain	faubourg Saint-Germain
Kraken	lc faubourg Saint-Gcrmain	Saint-Gcrmain	faubourg Saint-Gcrmain
Tess	le faubourg Saint-Germain	faubourg Saint-Germain	faubourg Saint-Germain
Tess fr	le faubourg Saint-Germain	faubourg Saint-Germain	faubourg Saint-Germain

Table – Alignement : Faux Positif ou Vrai Positif ?¹

* "Le chateau de Pinon", Dash , "Les trappeurs de l'Arkansas", Aimard.** "Souvenirs d'un vieux mélomane", Pontmartin.

Alignement manuel d'EN contaminées, OCR corrigés

	ELTeC Français ¹			
	OCR(Kraken)	Corr.	REN	REN corr.
MOBC	"[...] Paris"	"[...] Paris"	Paris Ici	Paris
MOMC	"[...] Grand-hlail"	"[...] Grand-Hail"	Grand-hlail	X
MOI	"[...] Mlorlincourtl"	"[...] Mlorlincourtl"	Mlorlincourtl	Mlorlincourtl
BOIC	"[...] Morlincourt "	"[...] Martincourt"	Morlincourt	Martincourt

	TGB ²			
	OCR	Corr.	REN	REN corr.
MOBC	"[...] Sainl-Cyr"	"[...] Saint-Cyr"	Sainl-Cyr	Saint-Cyr
MOMC	"[...] Nîme"	"[...] Mîme"	Nîme	Mîme
MOI	"[...] SaintAntoine"	"[...] SaintAntoine"	SaintAntoine	SaintAntoine
BOIC	"[...] Pinde"	"[...] Inde"	Pinde	Et du sommet du Inde

Table – Résultats de REN spaCy-ig. Types de formes : MOBC - mal océrisées bien corrigées ; MOMC - mal océrisées mal corrigées ; MOI - mal océrisées ignorées ; BOIC - bien océrisées indûment corrigées. Corr. - correction avec jampell, REN corr. - REN après la correction avec Jampell.

¹ *Mon Village*, Adam. "Le petit chose", Daudet.

² "Oeuvres poétiques de Boileau", Boileau. "Poésies de Valentin", Bourette. "Histoire des sociétés secrètes, politiques et religieuses", Zaccone.

Évaluer l'impact de la correction OCR sur la REN

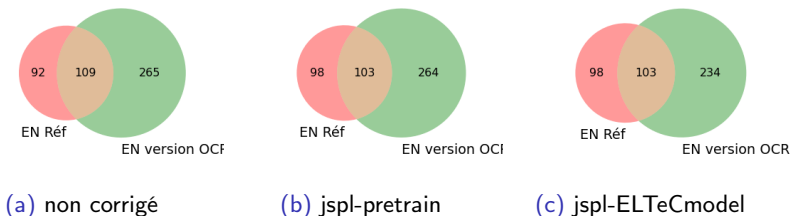


Figure – a), b), c) "Le petit chose", Daudet. Configuration : Kraken, spaCy 3.5.1, modèle lg

Évaluer l'impact de la correction OCR sur la REN

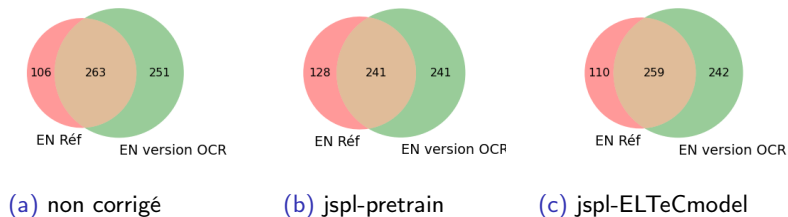
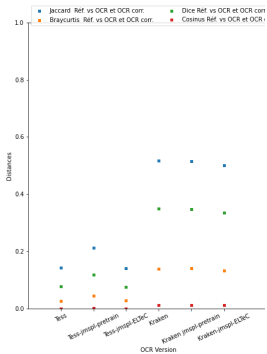
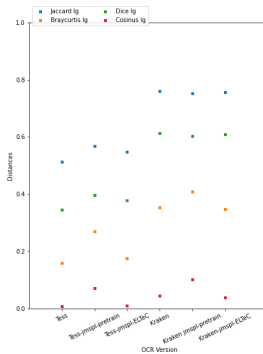


Figure – a), b), c) "*Le petit chose*", Daudet. Configuration : TesseractFra-Png, spaCy 3.5.1, modèle lg

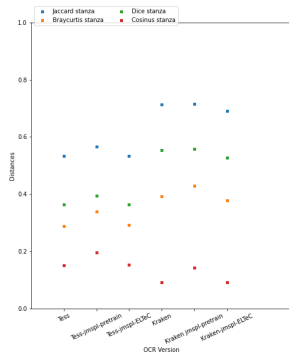
Évaluer l'impact de la correction OCR sur la REN



(a) texte



(b) spaCy_lg



(c) stanza

Figure – Distances Jaccard, Dice, Bray-Curtis, Cosinus, a), b), c) "Le petit chose", Daudet.

Évaluer l'impact de la correction OCR sur la REN

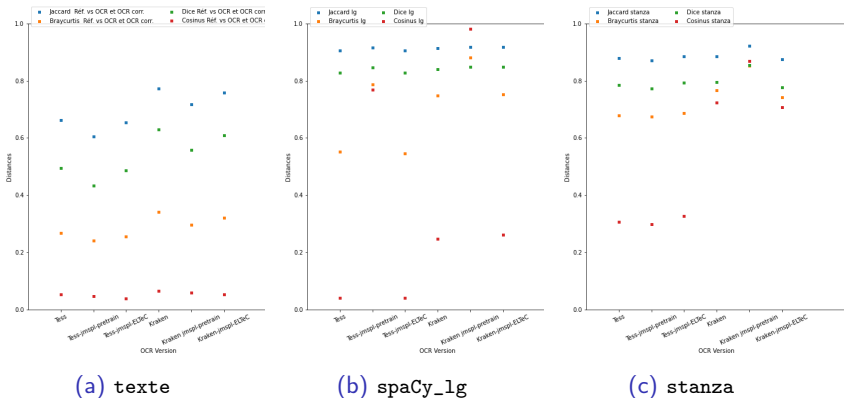


Figure – Distances Jaccard, Dice, Bray-Curtis, Cosinus, a), b), c) "Drame en 5 actes", Inconnu.

- OCR de meilleure qualité → Meilleurs résultats de REN ? → Moins de Faux Négatif/silence
- Correction de meilleure qualité → ? Amélioration de la REN ?

Nerval : un outil pour l'alignement des EN contaminées

Version	#Entités		Évaluation par Nerval			
	OCR	Réf.	Intersection	Précision	Rappel	F_1 mesure
Kraken	1122	944	566	0.504	0.761	0.607
Tess fr	860	944	646	0.751	0.868	0.805
Tess	920	944	597	0.649	0.802	0.718

Version	#Entités		Évaluation par Nerval			
	OCR	Réf.	Intersection	Précision	Rappel	F_1 mesure
Kraken + Jspl-ELTeC	1055	944	548	0.519	0.737	0.609 ↑
Tess fr + Jspl-ELTeC	838	944	621	0.741	0.835	0.785 ↓
Tess + Jspl-ELTeC	927	944	576	0.621	0.774	0.689 ↓

Version	#Entités		Évaluation par Nerval			
	OCR	Réf.	Intersection	Précision	Rappel	F_1 mesure
Kraken + Jspl-pt	1027	944	471	0.459	0.633	0.532 ↓
Tess fr + Jspl-pt	794	944	532	0.67	0.715	0.692 ↓
Tess + Jspl-pt	846	944	503	0.595	0.676	0.633 ↓

Table – Alignement des entités (spaCy_{lg} format IOB) de la référence et des versions OCR corrigées avec le modèle français préentraîné de Jampell et un modèle entraîné sur une partie du corpus ELTeC, Nerval⁴, "Le petit chose", Daudet, 1868.

4. <https://gitlab.com/tekliia/ner/nerval>

- **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
 - **EN contaminées = Variations orthographiques** d'un toponyme
- Usages de **différentes métriques** : évaluation fine/nuancée
- Méthodes d'alignements et évaluation : évaluation Nerval

- **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
 - **EN contaminées = Variations orthographiques** d'un toponyme
- Usages de **différentes métriques** : évaluation fine/nuancée
- Méthodes d'alignements et évaluation : évaluation Nerval
- **La correction automatique n'est pas automatiquement un gain** :
 - l'outil produit un nombre conséquent des sur-corrections
 - l'outil ne corrige pas toutes les contaminations/interférences OCR
 - La correction des EN a moins d'impact que la correction du contexte(?)
- **Typologie des erreurs de corrections automatiques d'OCR**